

# Distributed Systems

CS425/ECE428

*Instructor: Radhika Mittal*

*Acknowledgements for some of materials: Indy Gupta and Nikita Borisov*

# Logistics

- MP0 is due today at 11:59pm.
- Reminder to share your name when you speak up in class.
- Feb 7 lecture was not recorded! Please see my Campuswire post for links to last year's videos covering the same topic.

# Today's agenda

- **Global State (contd.)**
  - Chapter 14.5
  
- **Multicast**
  - Chapter 15.4

# Today's agenda

- **Global State (contd.)**
  - Chapter 14.5
  
- **Multicast**
  - Chapter 15.4

# Global Snapshot Summary

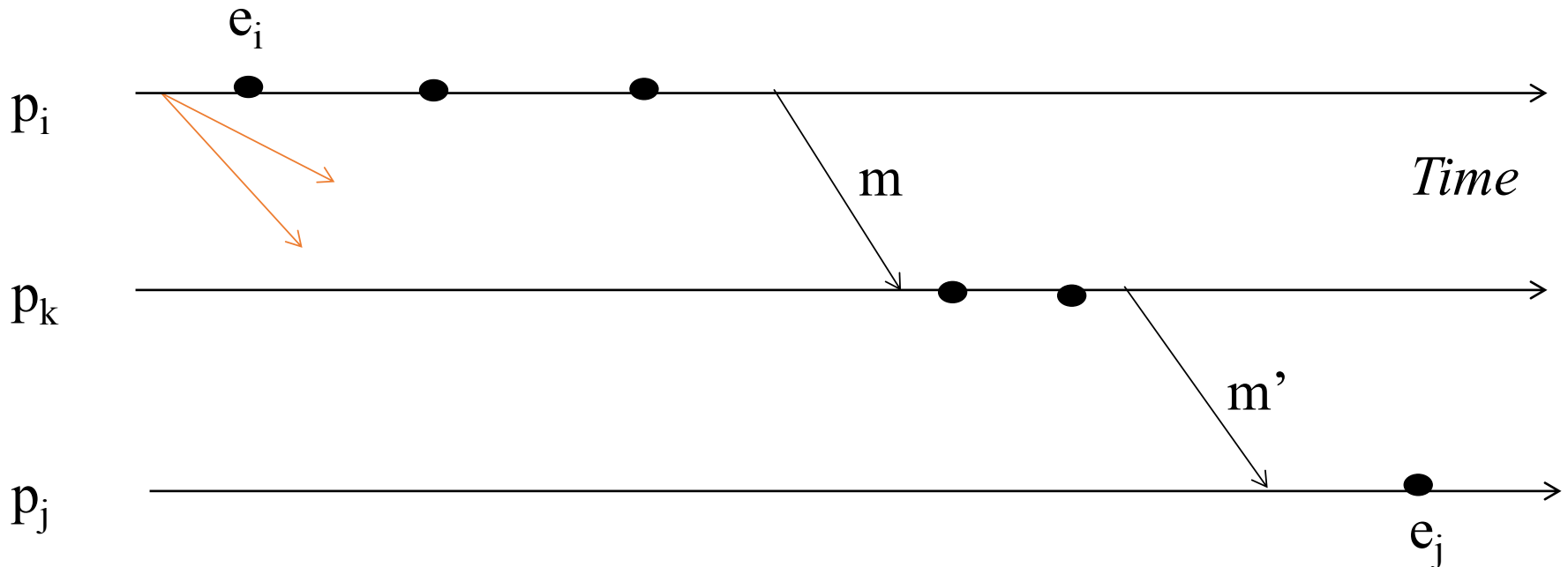
- The ability to calculate global snapshots in a distributed system is very important.
- But don't want to interrupt running distributed application.
- Chandy-Lamport algorithm calculates global snapshot.
- Obeys causality (creates a consistent cut).

# Chandy-Lamport Algorithm: Properties

- Any run of the Chandy-Lamport Global Snapshot algorithm creates a consistent cut.
- Let  $e_i$  and  $e_j$  be events occurring at  $p_i$  and  $p_j$ , respectively such that
  - $e_i \rightarrow e_j$  ( $e_i$  happens before  $e_j$ )
- The snapshot algorithm ensures that
  - if  $e_j$  is in the cut then  $e_i$  is also in the cut.
- That is: if  $e_j \rightarrow \langle p_j \text{ records its state} \rangle$ , then it must be true that  $e_i \rightarrow \langle p_i \text{ records its state} \rangle$ .

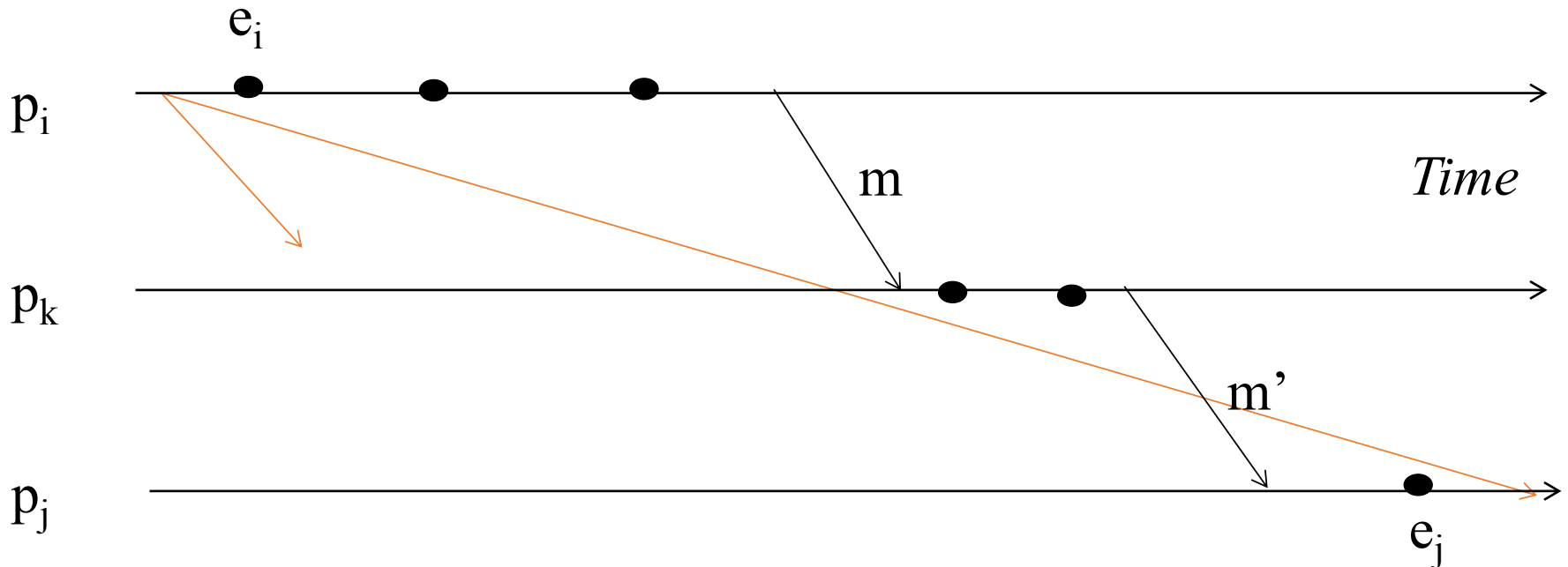
# Chandy-Lamport Algorithm: Properties

- If  $e_j \rightarrow \langle p_j \text{ records its state} \rangle$ , then it must be true that  $e_i \rightarrow \langle p_i \text{ records its state} \rangle$ .
- By contradiction, suppose  $e_j \rightarrow \langle p_j \text{ records its state} \rangle$ , and  $\langle p_i \text{ records its state} \rangle \rightarrow e_i$ .



# Chandy-Lamport Algorithm: Properties

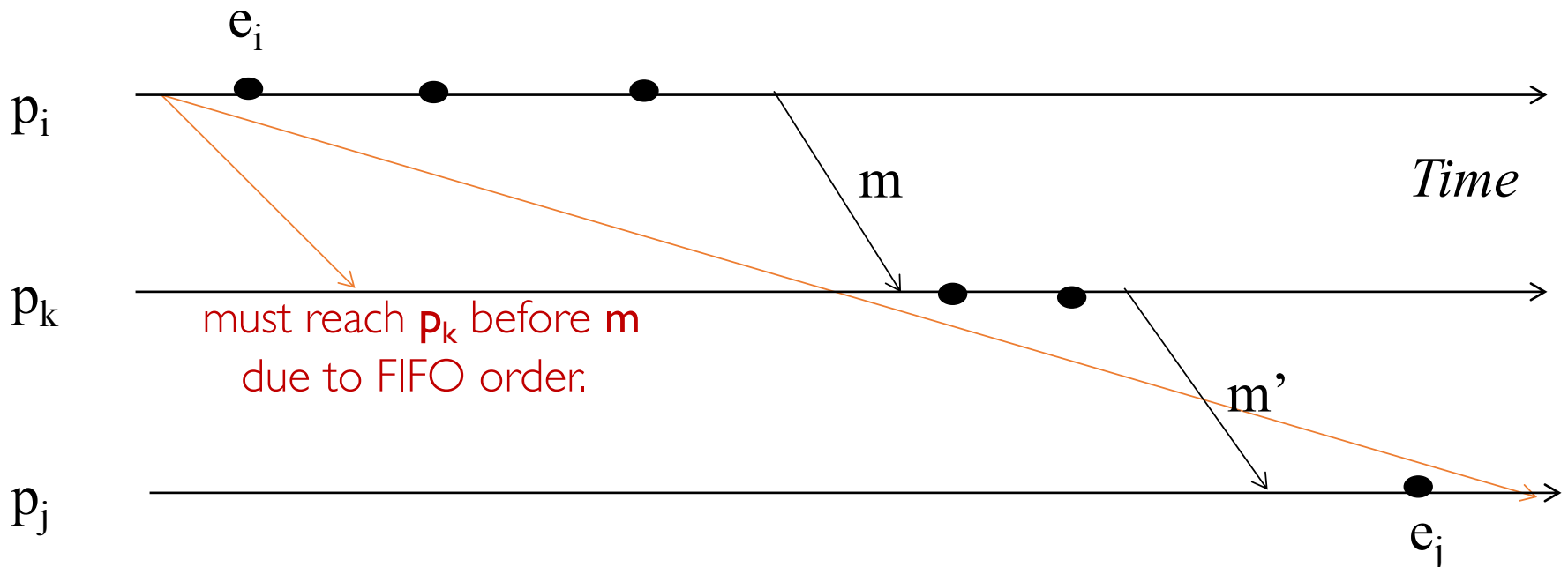
- If  $e_j \rightarrow \langle p_j \text{ records its state} \rangle$ , then it must be true that  $e_i \rightarrow \langle p_i \text{ records its state} \rangle$ .
- By contradiction, suppose  $e_j \rightarrow \langle p_j \text{ records its state} \rangle$ , and  $\langle p_i \text{ records its state} \rangle \rightarrow e_i$ .





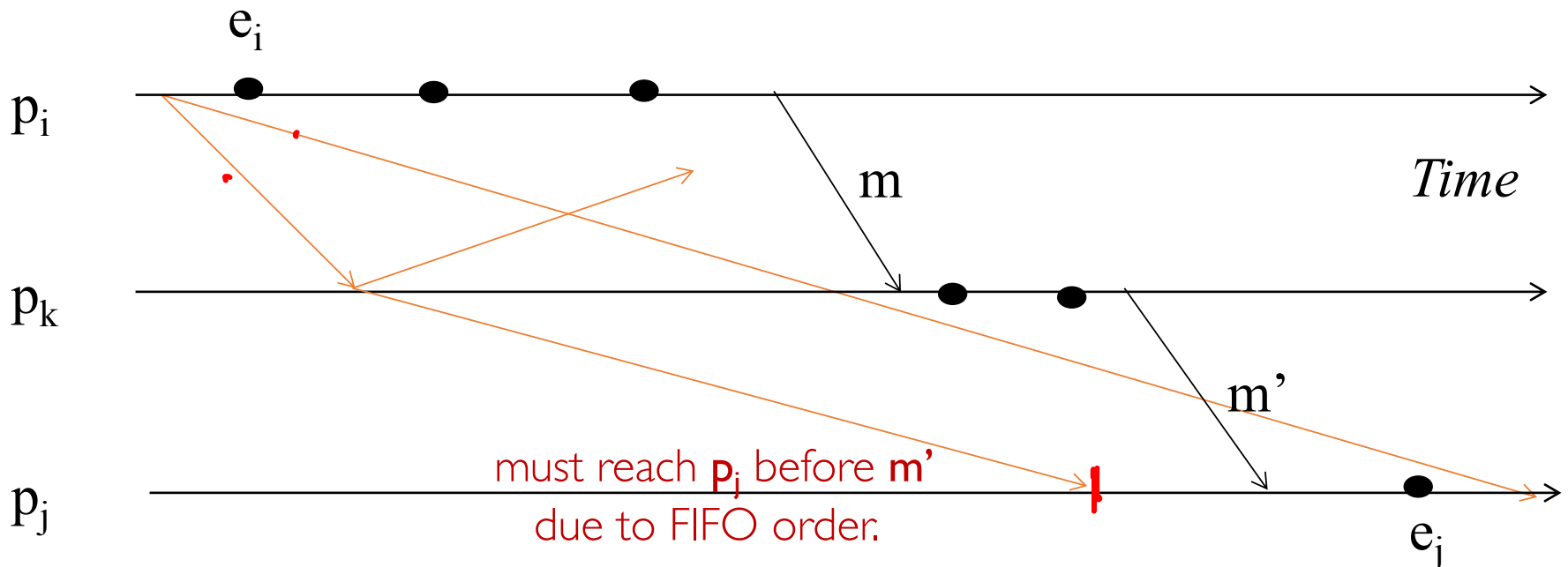
# Chandy-Lamport Algorithm: Properties

- If  $e_j \rightarrow \langle p_j \text{ records its state} \rangle$ , then it must be true that  $e_i \rightarrow \langle p_i \text{ records its state} \rangle$ .
- By contradiction, suppose  $e_j \rightarrow \langle p_j \text{ records its state} \rangle$ , and  $\langle p_i \text{ records its state} \rangle \rightarrow e_i$ .



# Chandy-Lamport Algorithm: Properties

- If  $e_j \rightarrow \langle p_j \text{ records its state} \rangle$ , then it must be true that  $e_i \rightarrow \langle p_i \text{ records its state} \rangle$ .
- By contradiction, suppose  $e_j \rightarrow \langle p_j \text{ records its state} \rangle$ , and  $\langle p_i \text{ records its state} \rangle \rightarrow e_i$ .



# Chandy-Lamport Algorithm: Properties

- If  $e_j \rightarrow \langle p_j \text{ records its state} \rangle$ , then it must be true that  $e_i \rightarrow \langle p_i \text{ records its state} \rangle$ .
- By contradiction, suppose  $e_j \rightarrow \langle p_j \text{ records its state} \rangle$ , and  $\langle p_i \text{ records its state} \rangle \rightarrow e_i$ .
- Consider the path of app messages (through other processes) that go from  $e_i$  to  $e_j$ .
- Due to FIFO ordering, markers on each link in above path will precede regular app messages.
- Thus, since  $\langle p_i \text{ records its state} \rangle \rightarrow e_i$ , it must be true that  $p_j$  received a marker before  $e_j$ .
- Thus  $e_j$  is not in the cut  $\Rightarrow$  contradiction.

# Why capture global snapshots?

- *Checkpointing* the system state.
- Reasoning about unreferenced objects (for garbage collection).
- Distributed debugging.
- Can be used to detect **global properties**.
  - **Safety vs. Liveness**

# Revisions: notations and definitions

- For a process  $p_i$ , where events  $e_i^1, \dots$  occur:

$$\text{history}(p_i) = h_i = \langle e_i^1, \dots \rangle$$

$$\text{prefix history}(p_i^k) = h_i^k = \langle e_i^1, \dots, e_i^k \rangle$$

$s_i^k$ :  $p_i$ 's state immediately after  $k^{\text{th}}$  event.

- For a set of processes  $\langle p_1, p_2, p_3, \dots, p_n \rangle$ :

$$\text{global history: } H = \cup_i (h_i)$$

$$\text{a cut } C \subseteq H = h_1^{c_1} \cup h_2^{c_2} \cup \dots \cup h_n^{c_n}$$

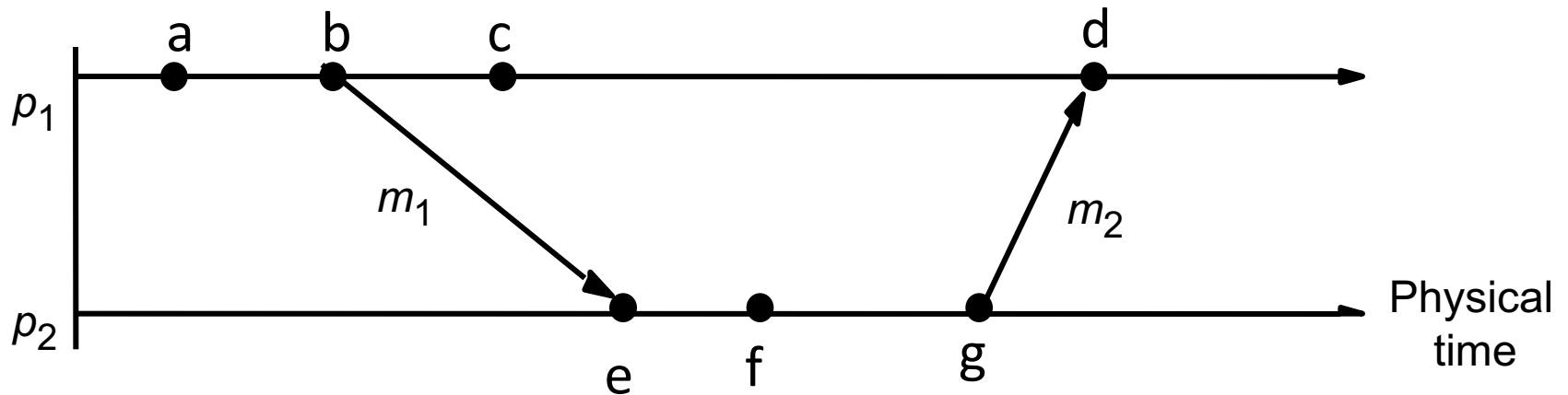
$$\text{the frontier of } C = \{e_i^{c_i}, i = 1, 2, \dots, n\}$$

$$\text{global state } S \text{ that corresponds to cut } C = \cup_i (s_i^{c_i})$$

# More notations and definitions

- A **run** is a total ordering of events in  $H$  that is consistent with each  $h_i$ 's ordering.
- A **linearization** is a run consistent with happens-before ( $\rightarrow$ ) relation in  $H$ .

# Example



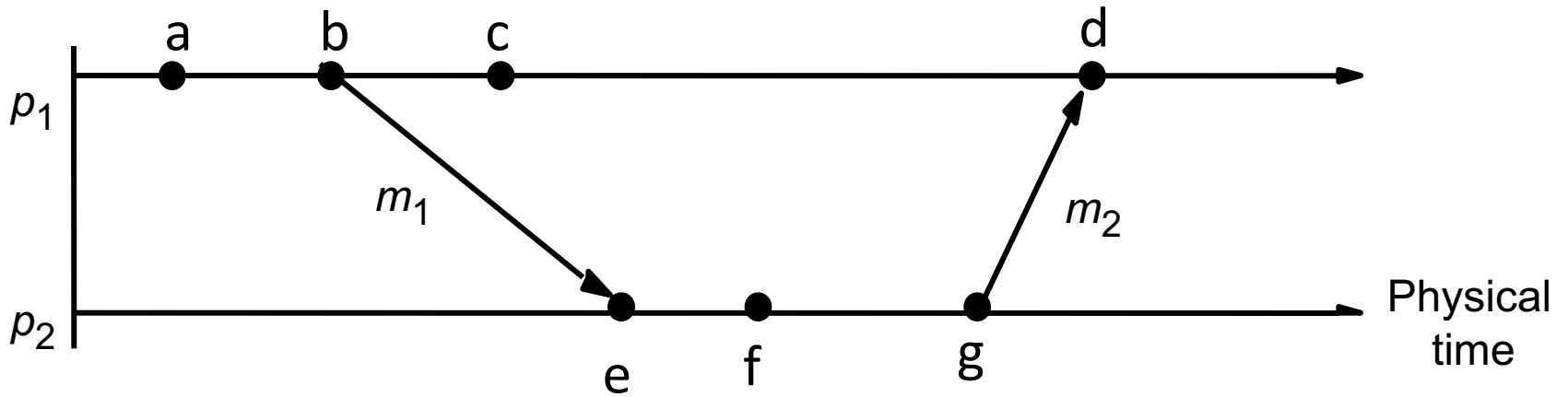
Order at  $p_1$ :  $\langle a, b, c, d \rangle$       Order at  $p_2$ :  $\langle e, f, g \rangle$

Causal order across  $p_1$  and  $p_2$ :  $\langle b, e \rangle, \langle g, d \rangle$

Run:  $\langle a, b, c, d, e, f, g \rangle \rightarrow$

Linearization:  $\langle a, b, c, e, f, g, d \rangle$

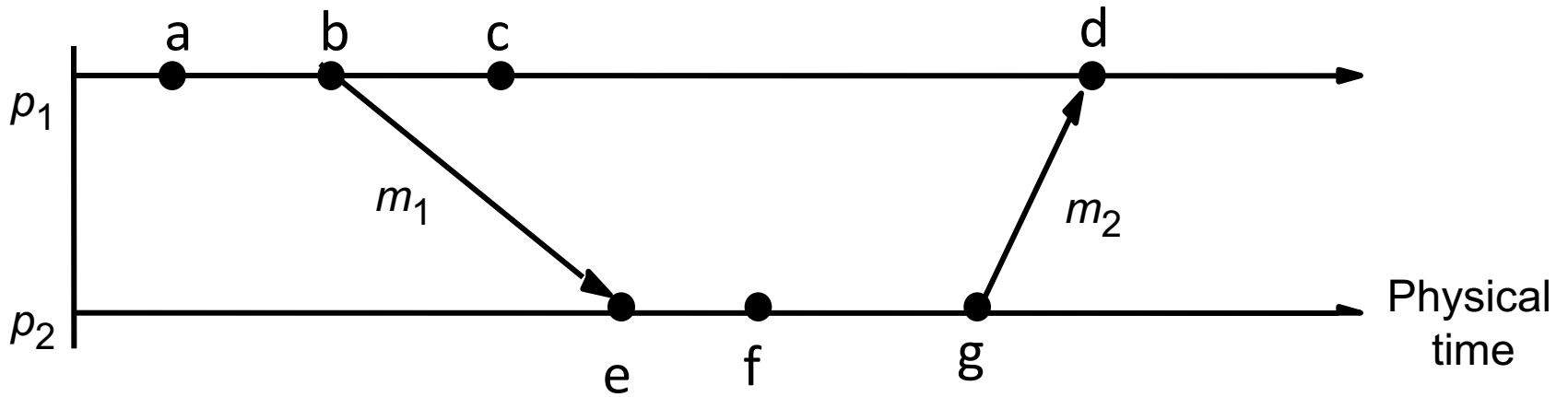
# Example



$\langle a, b, e, f, c, g, d \rangle$ :

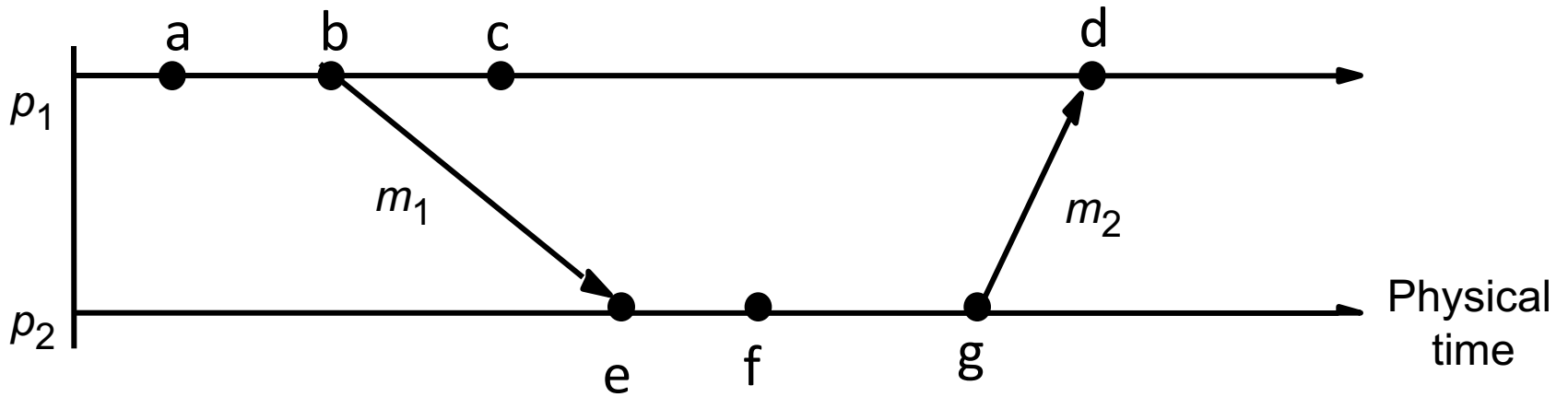


# Example



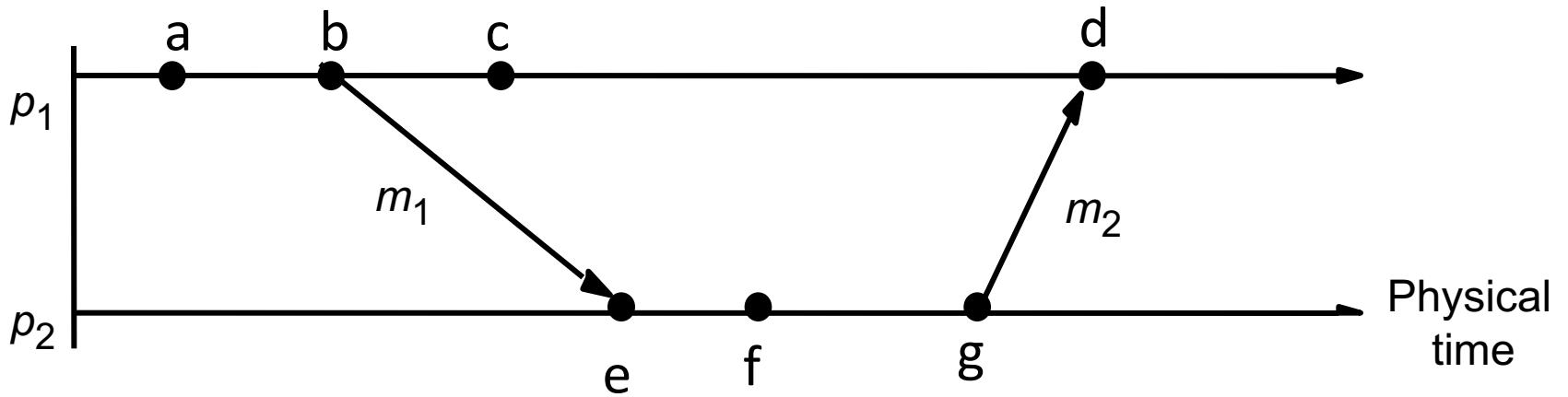
$\langle a, b, e, f, c, g, d \rangle$ : Linearization

# Example



$\langle a, f, e, b, c, g, d \rangle$ :

# Example

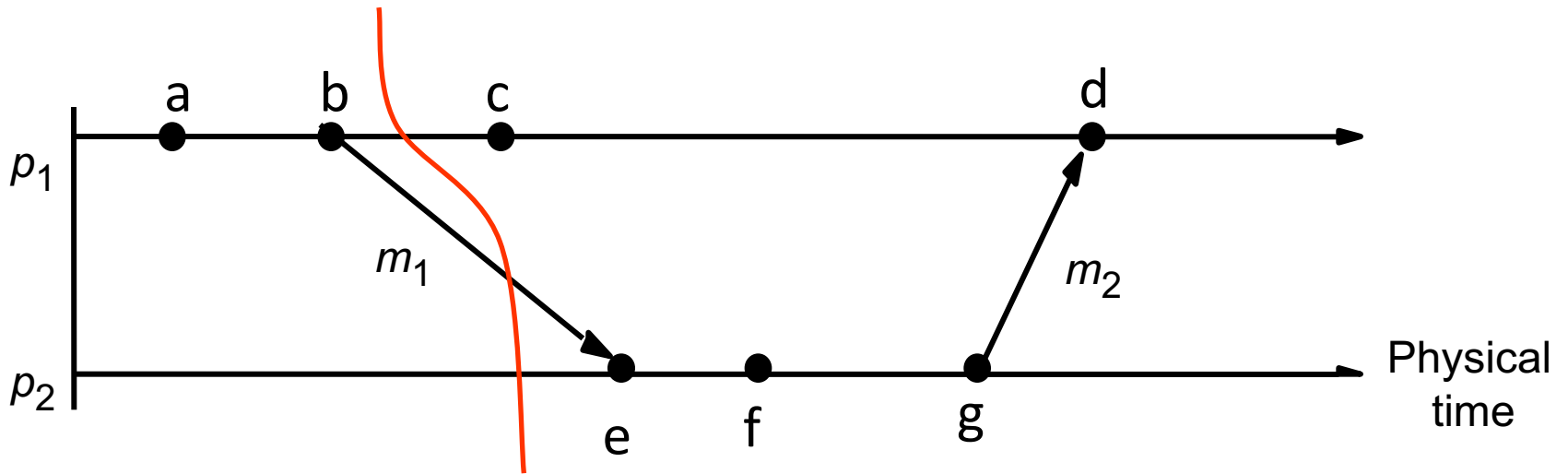


< a, f, e, b, c, g, d >: **Not even a run**

# More notations and definitions

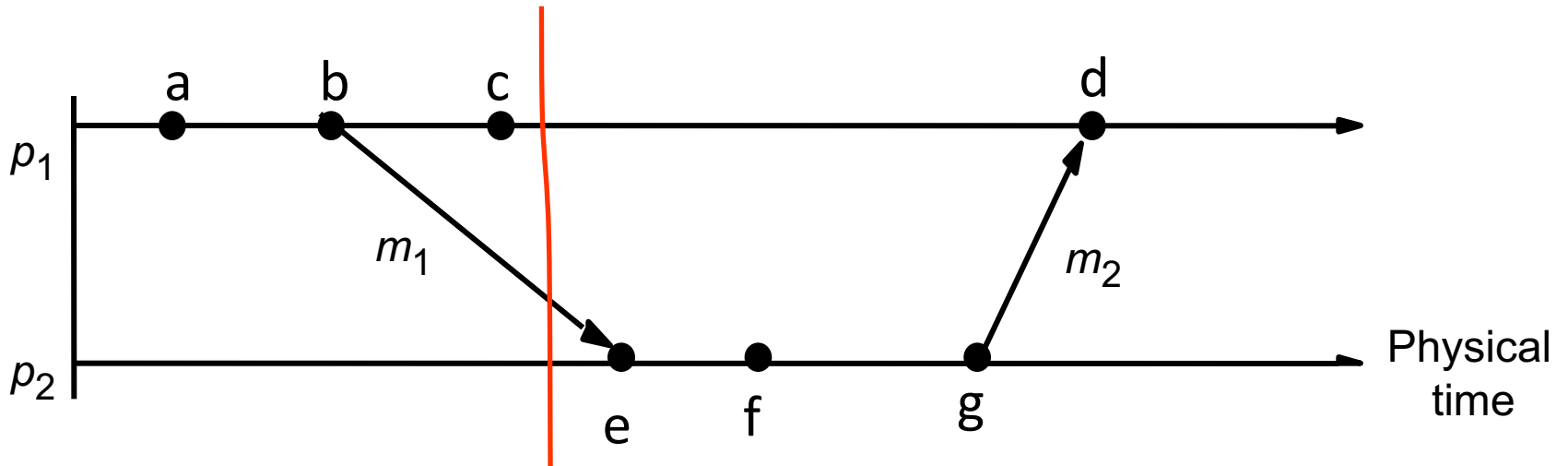
- A **run** is a total ordering of events in  $H$  that is consistent with each  $h_i$ 's ordering.
- A **linearization** is a run consistent with happens-before ( $\rightarrow$ ) relation in  $H$ .
- Linearizations pass through consistent global states.

# Example



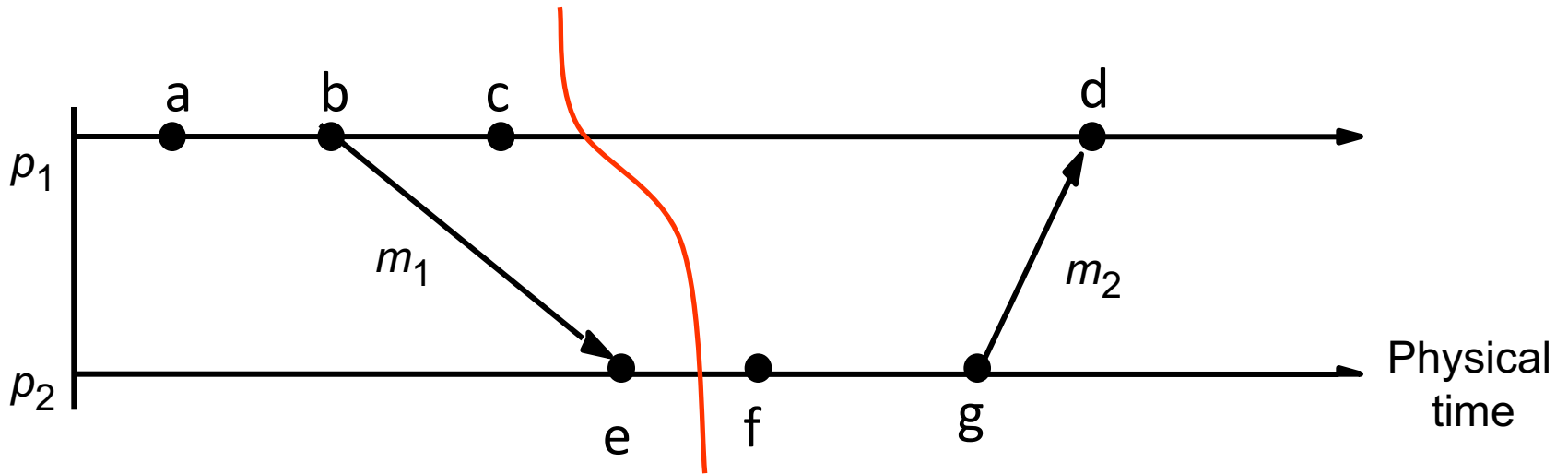
Linearization:  $\langle \underline{a, b}, |c, e, f, g, d \rangle$

# Example



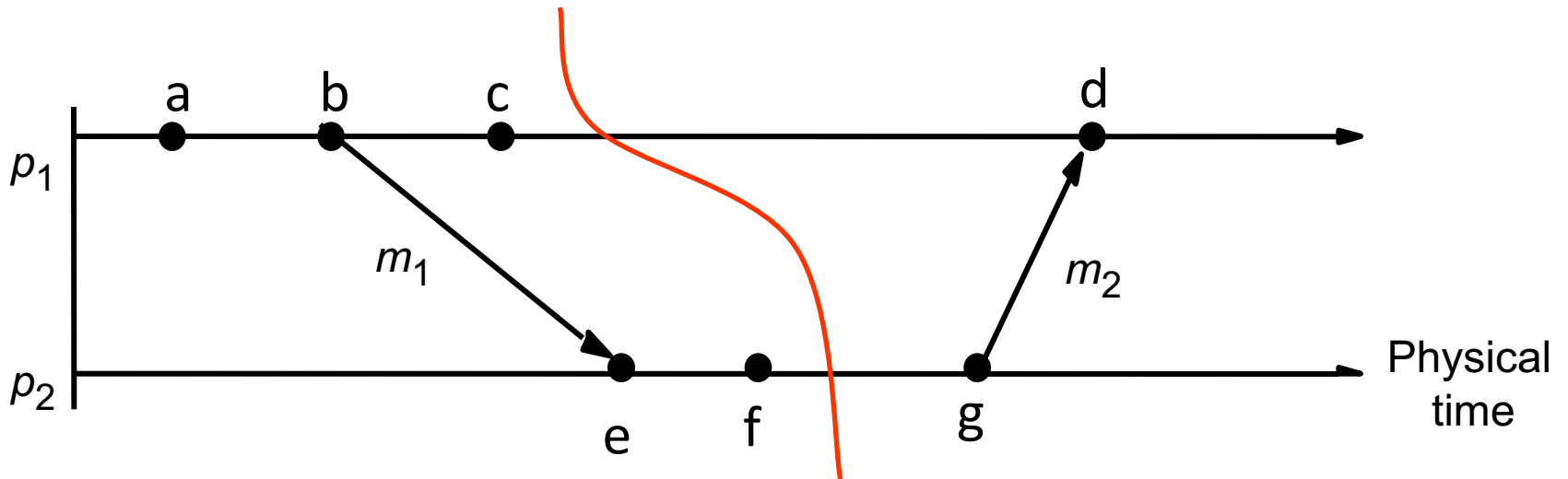
Linearization:  $\langle a, b, c, |e, f, g, d \rangle$

# Example



Linearization:  $\langle a, b, c, e, f, g, d \rangle$

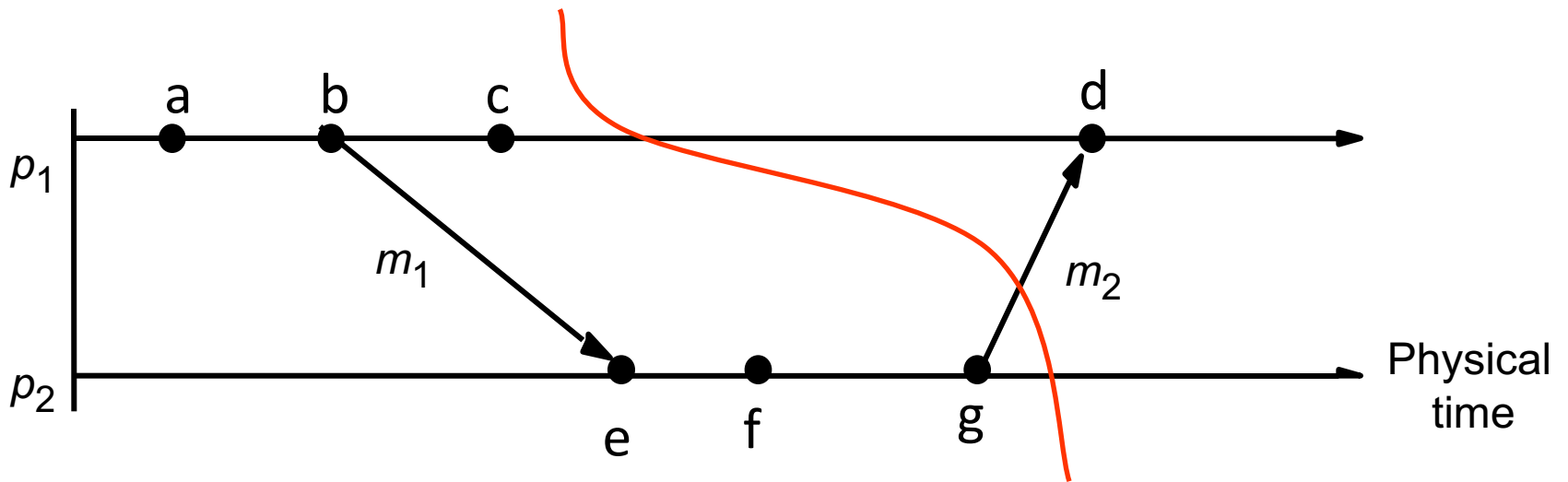
# Example



Linearization:  $\langle a, b, c, e, f, \mathbf{g}, d \rangle$

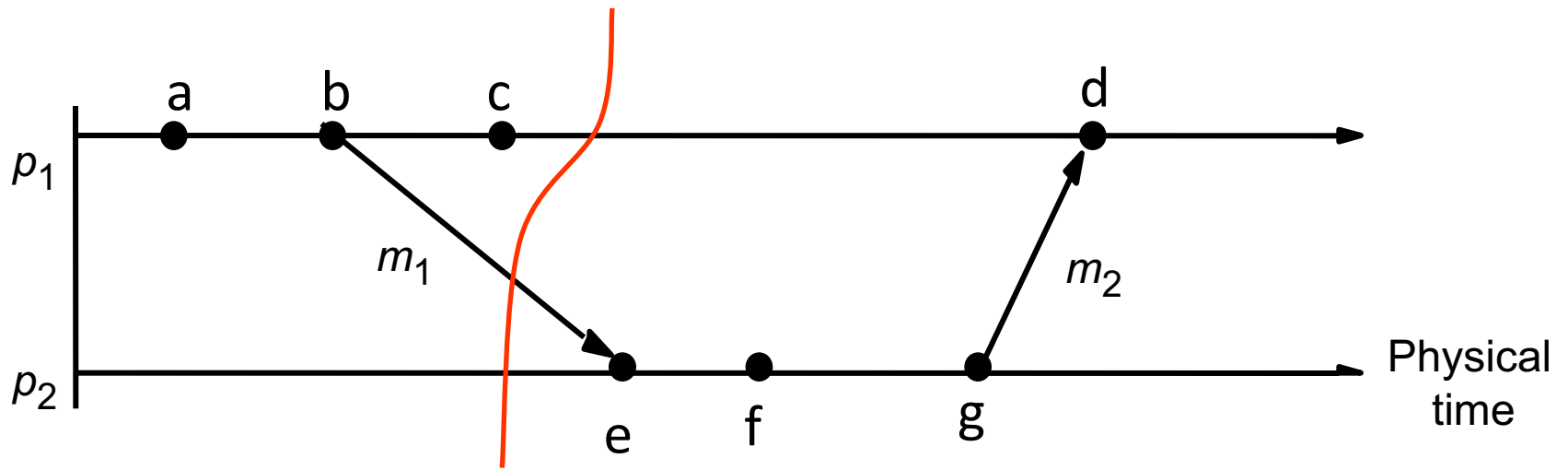


# Example



Linearization:  $\langle a, b, c, e, f, g \mid d \rangle$

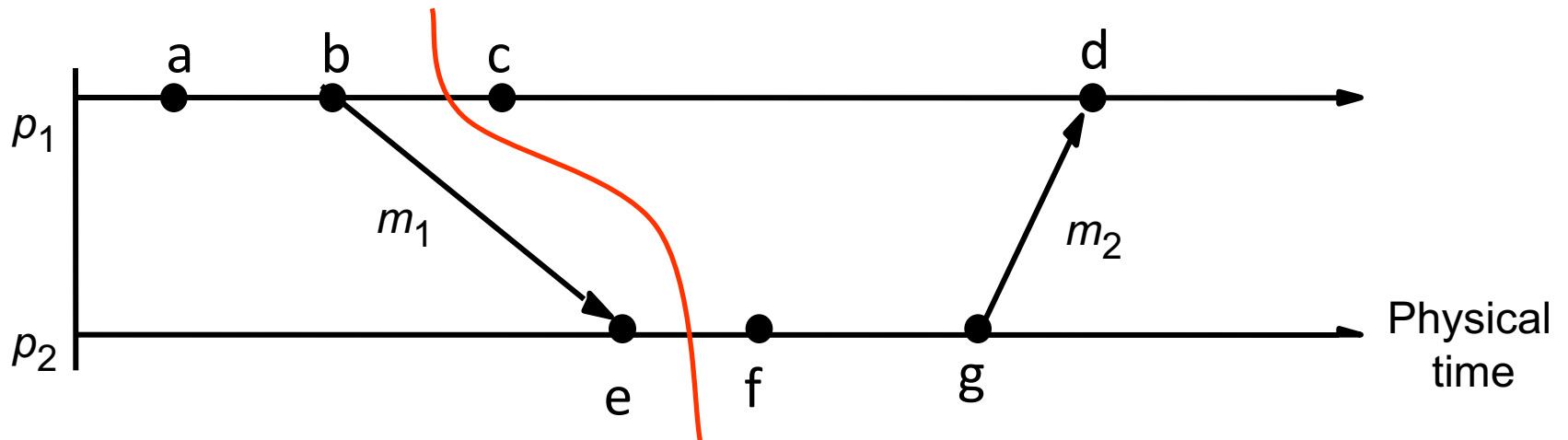
# Example



Linearization:  $\langle a, b, c, e, f, g, d \rangle$

Linearization:  $\langle a, b, e, c, f, g, d \rangle$

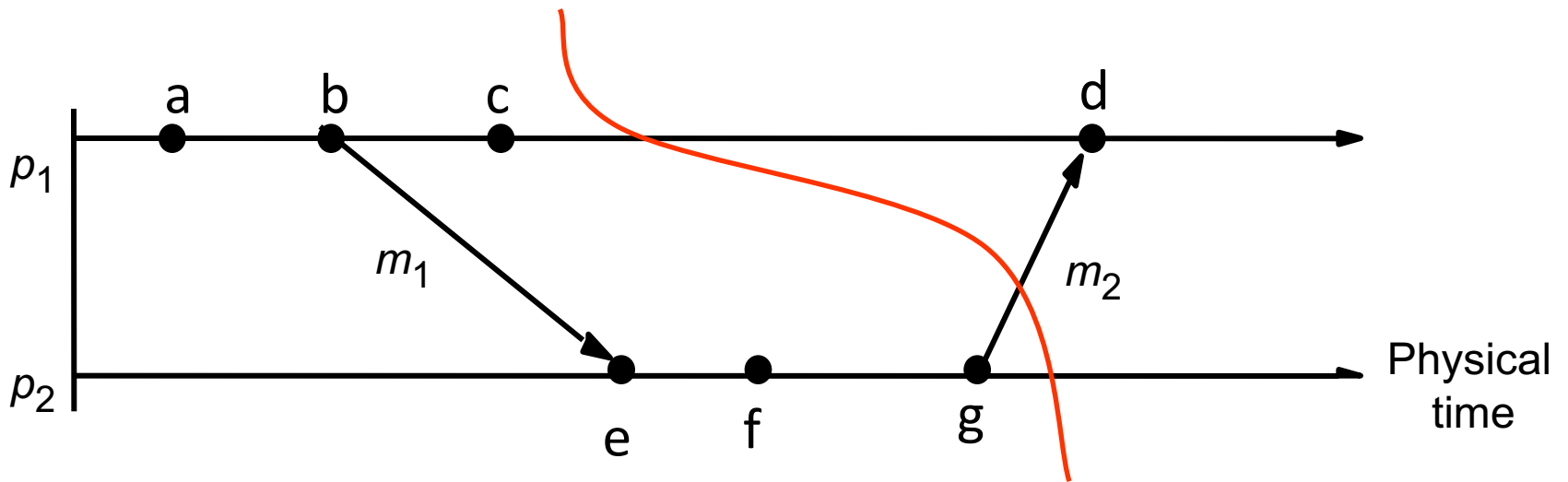
# Example



Linearization:  $\langle a, b, c, e, f, g, d \rangle$

Linearization:  $\langle a, b, e, c, f, g, d \rangle$

# Example



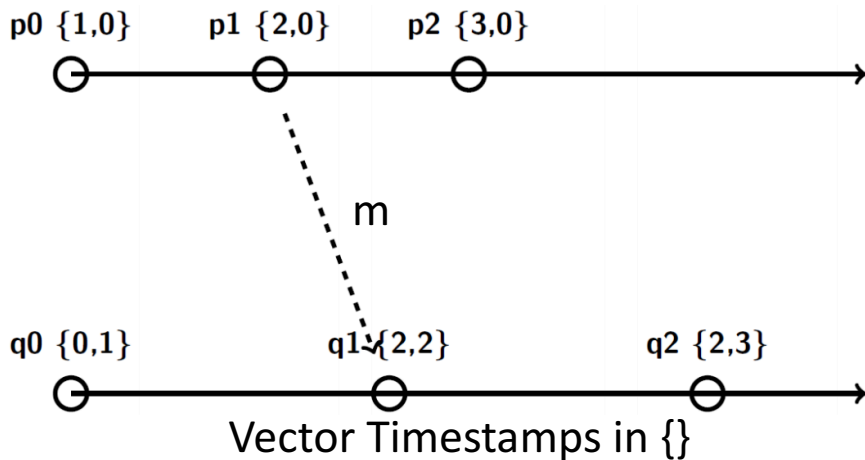
Linearization:  $\langle a, b, c, e, f, g \mid d \rangle$

Linearization:  $\langle a, b, e, c, f, g \mid d \rangle$

# More notations and definitions

- A **run** is a total ordering of events in  $H$  that is consistent with each  $h_i$ 's ordering.
- A **linearization** is a run consistent with happens-before ( $\rightarrow$ ) relation in  $H$ .
- Linearizations pass through consistent global states.
- A global state  $S_k$  is reachable from global state  $S_i$ , if there is a linearization that passes through  $S_i$  and then through  $S_k$ .
- The distributed system evolves as a series of transitions between global states  $S_0, S_1, \dots$

# State Transitions: Example



- Causal order:
  - $p_0 \rightarrow p_1 \rightarrow p_2$
  - $q_0 \rightarrow q_1 \rightarrow q_2$
  - $p_0 \rightarrow p_1 \rightarrow q_1 \rightarrow q_2$

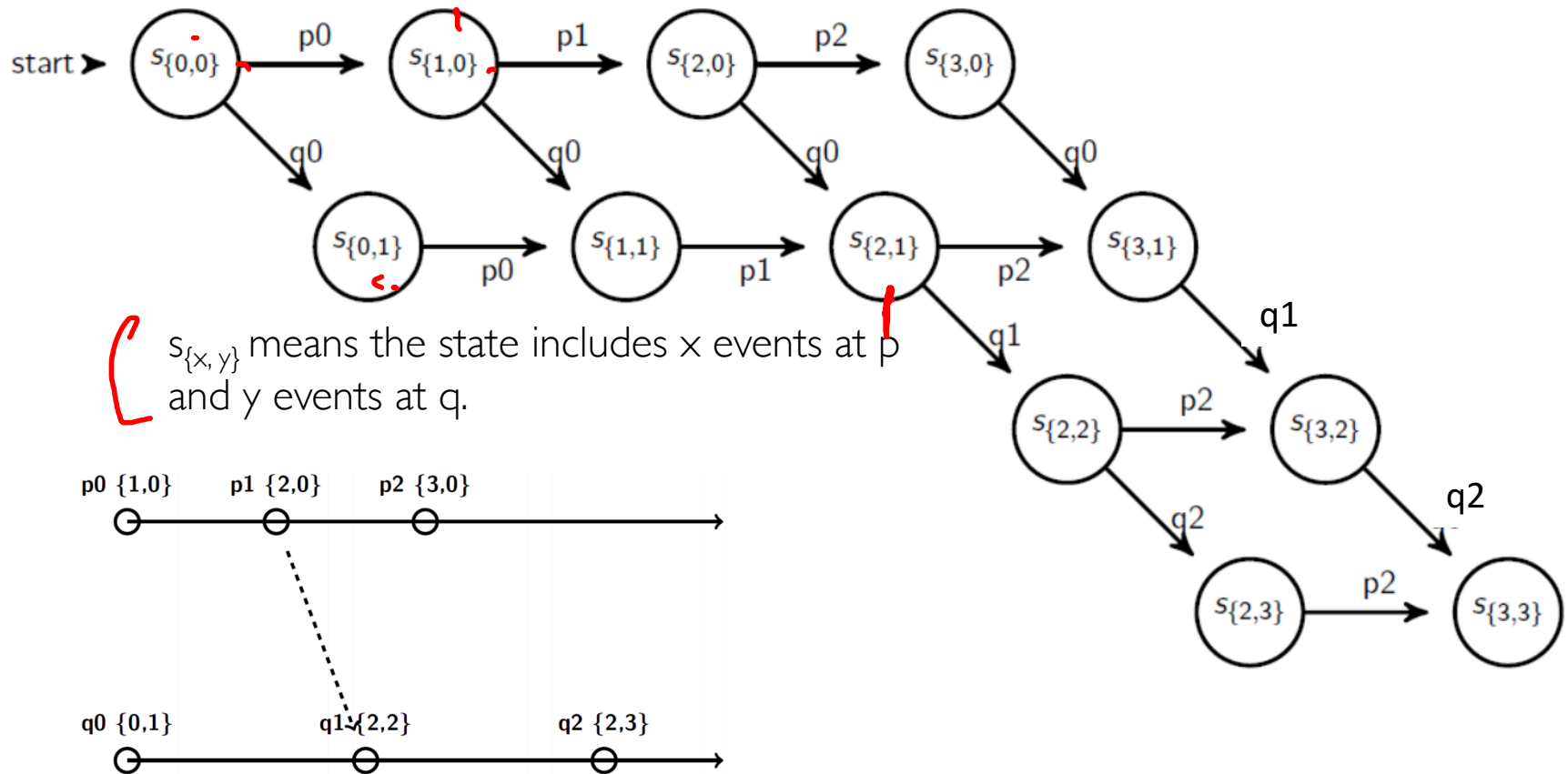
- Concurrent:
  - $p_0 \parallel q_0$
  - $p_1 \parallel q_0$
  - $p_2 \parallel q_0, p_2 \parallel q_1, p_2 \parallel q_2$

Many linearizations:

- $\langle p_0, p_1, p_2, q_0, q_1, q_2 \rangle$
- $\langle p_0, q_0, p_1, q_1, p_2, q_2 \rangle$
- $\langle q_0, p_0, p_1, q_1, p_2, q_2 \rangle$
- $\langle q_0, p_0, p_1, p_2, q_1, q_2 \rangle$
- .....

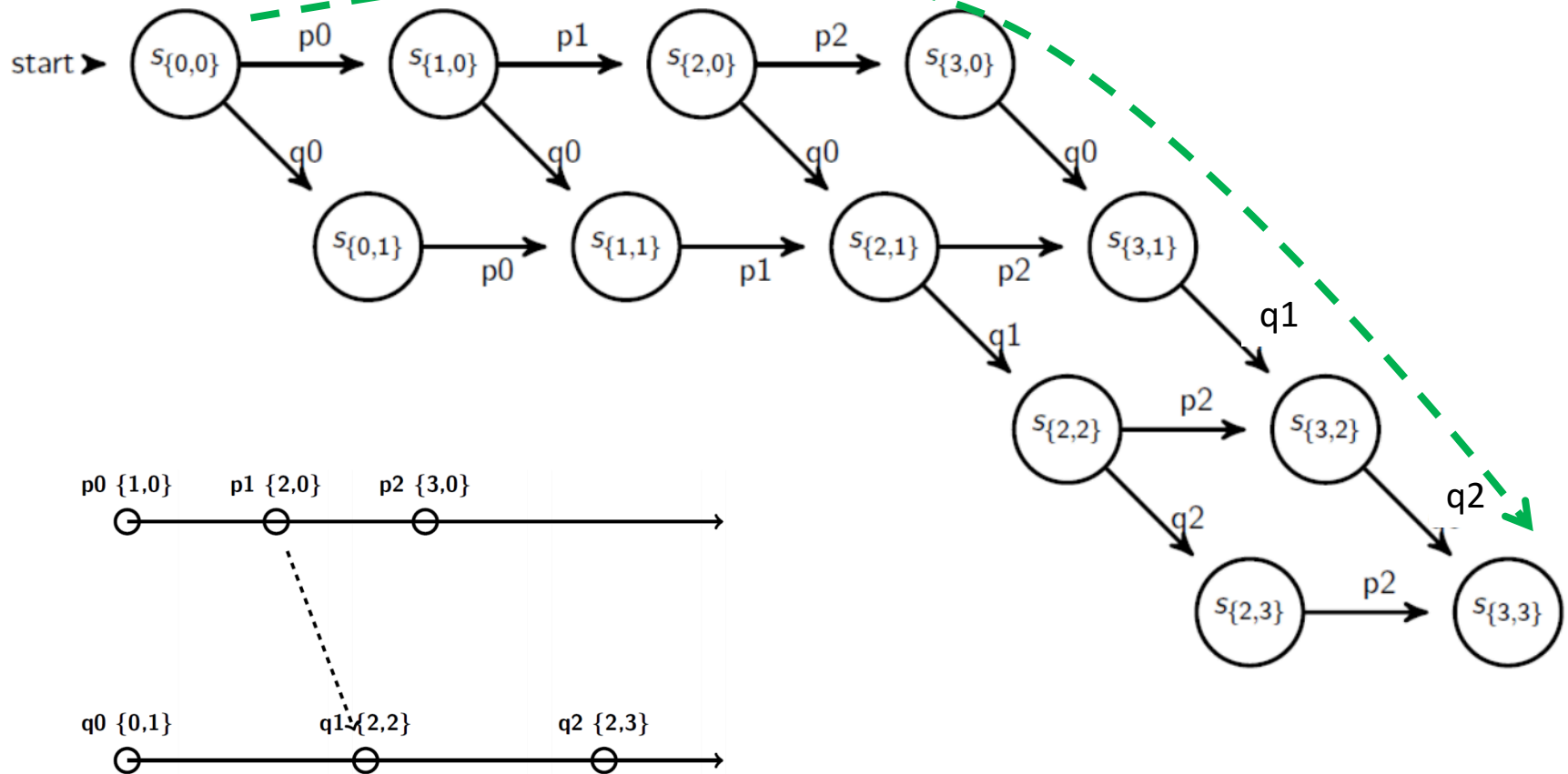
# State Transitions: Example

**Execution Lattice.** Each path represents a linearization.



# State Transitions: Example

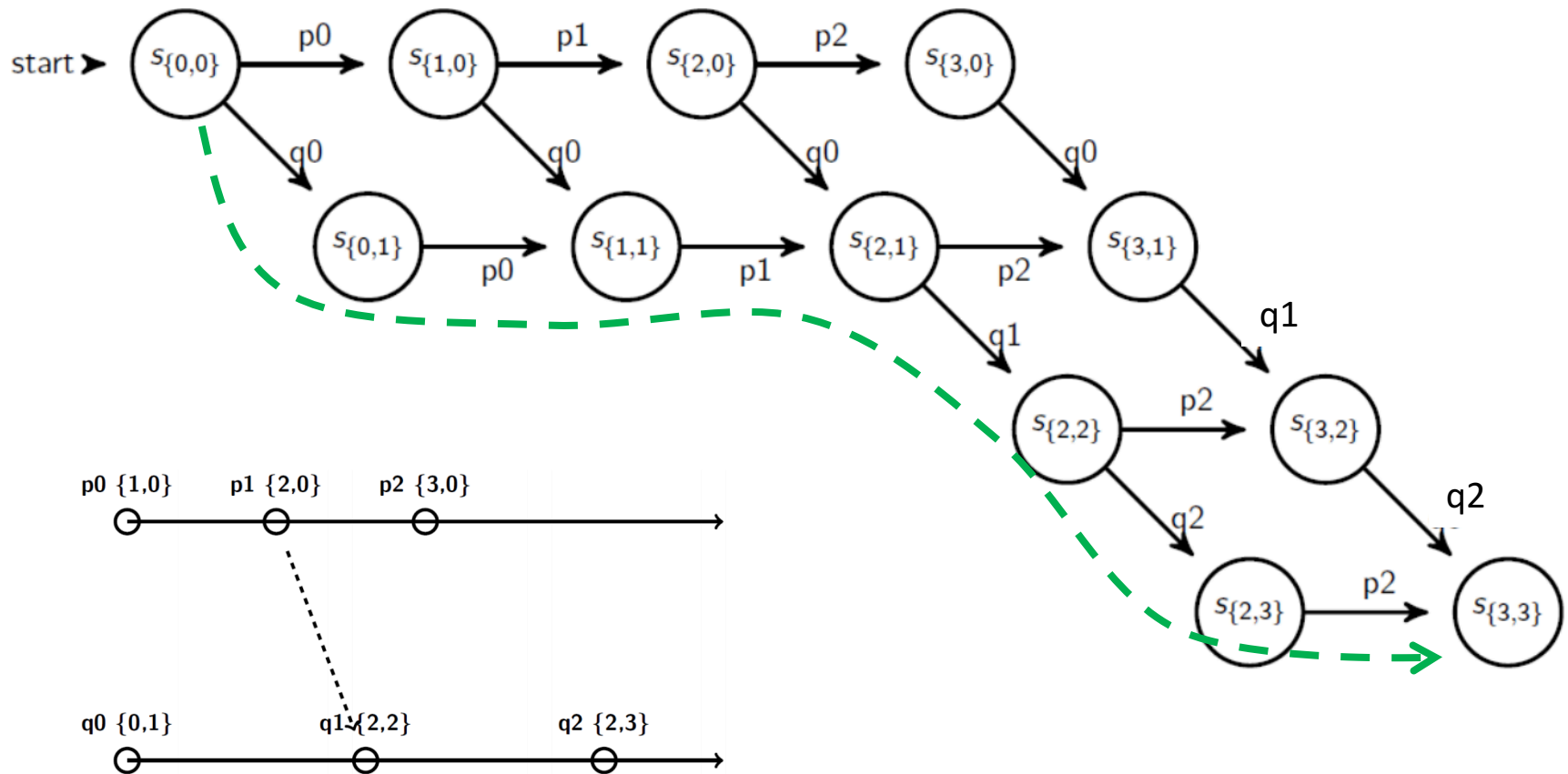
**Execution Lattice.** Each path represents a linearization.





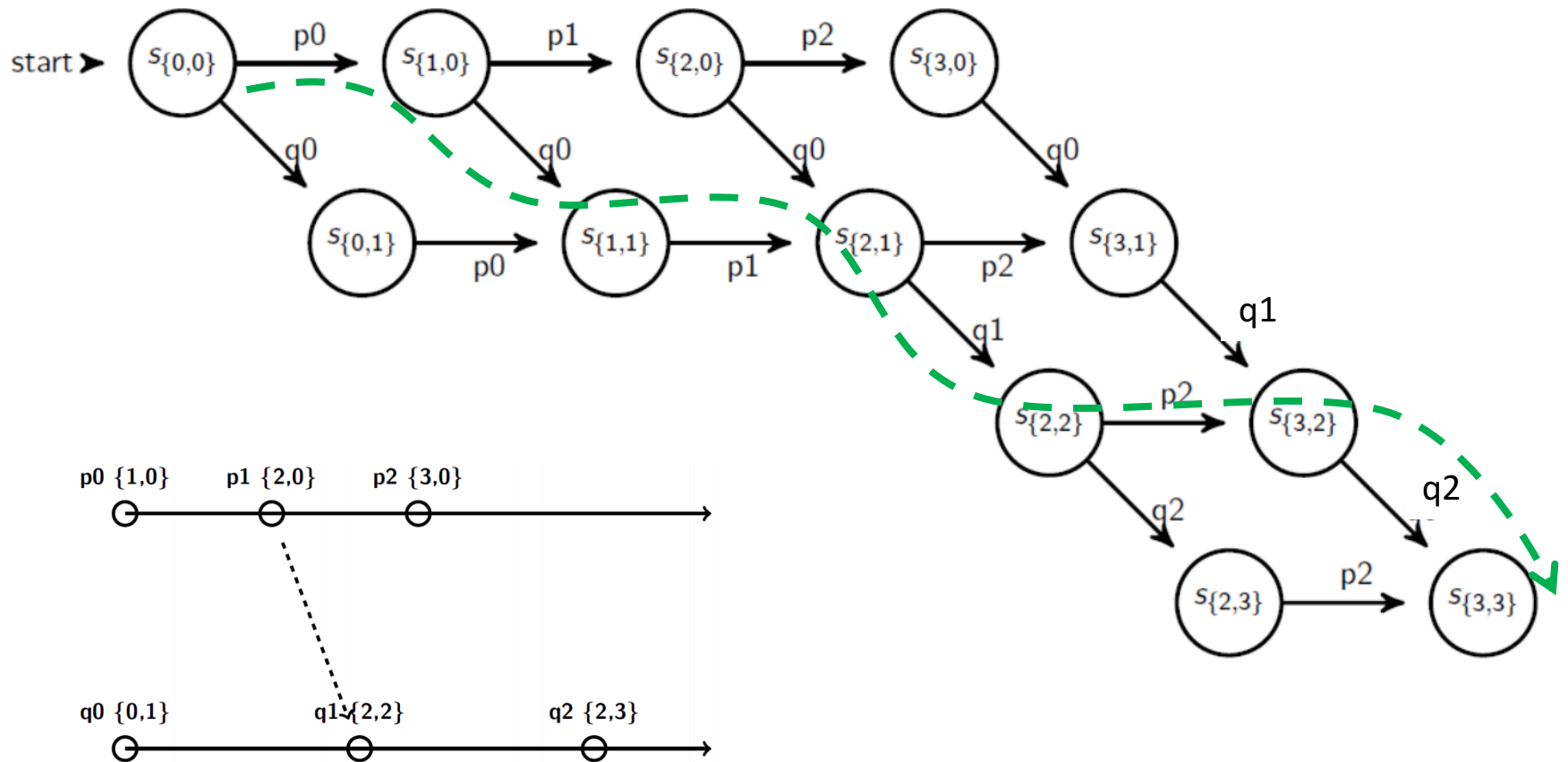
# State Transitions: Example

**Execution Lattice.** Each path represents a linearization.



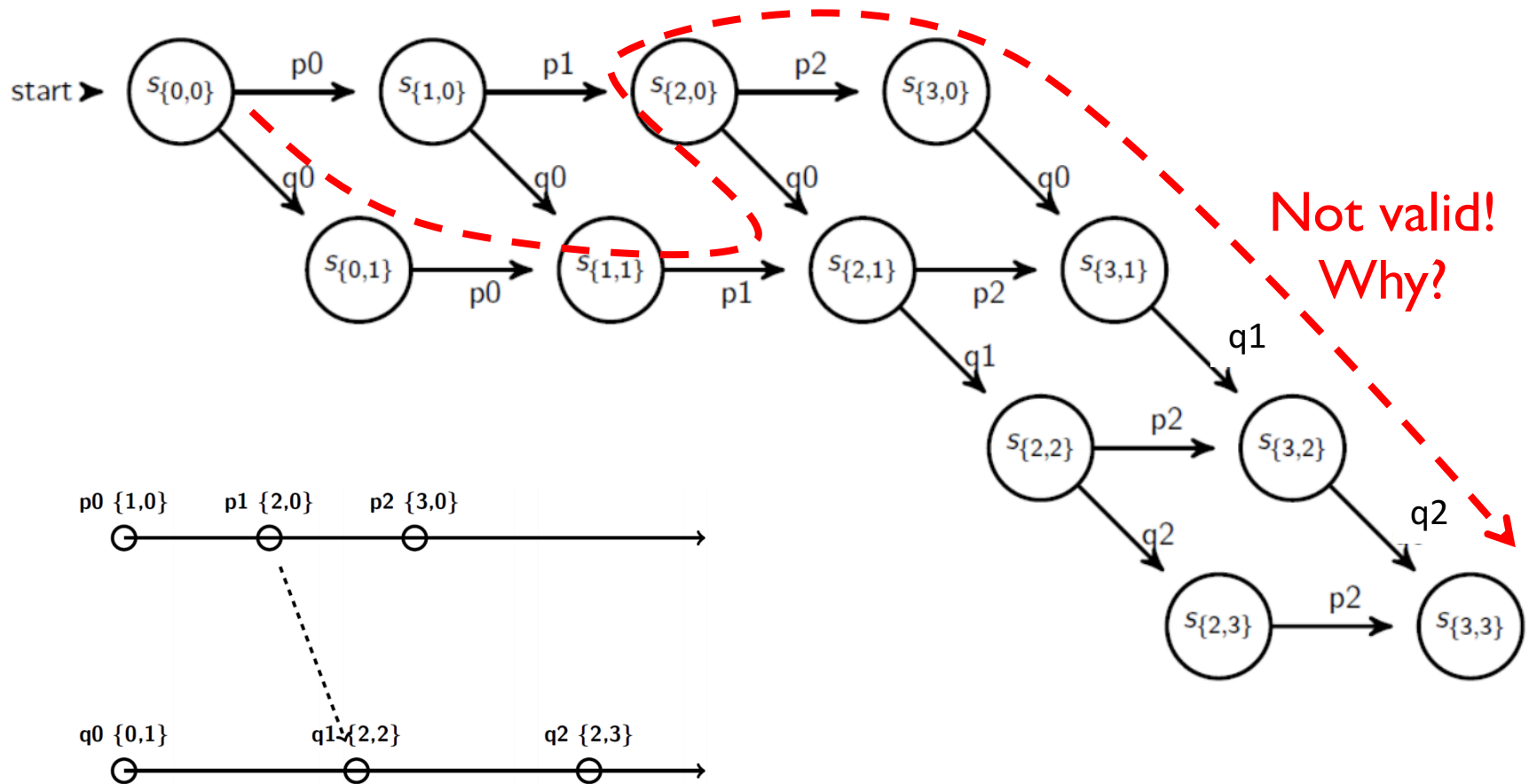
# State Transitions: Example

**Execution Lattice.** Each path represents a linearization.

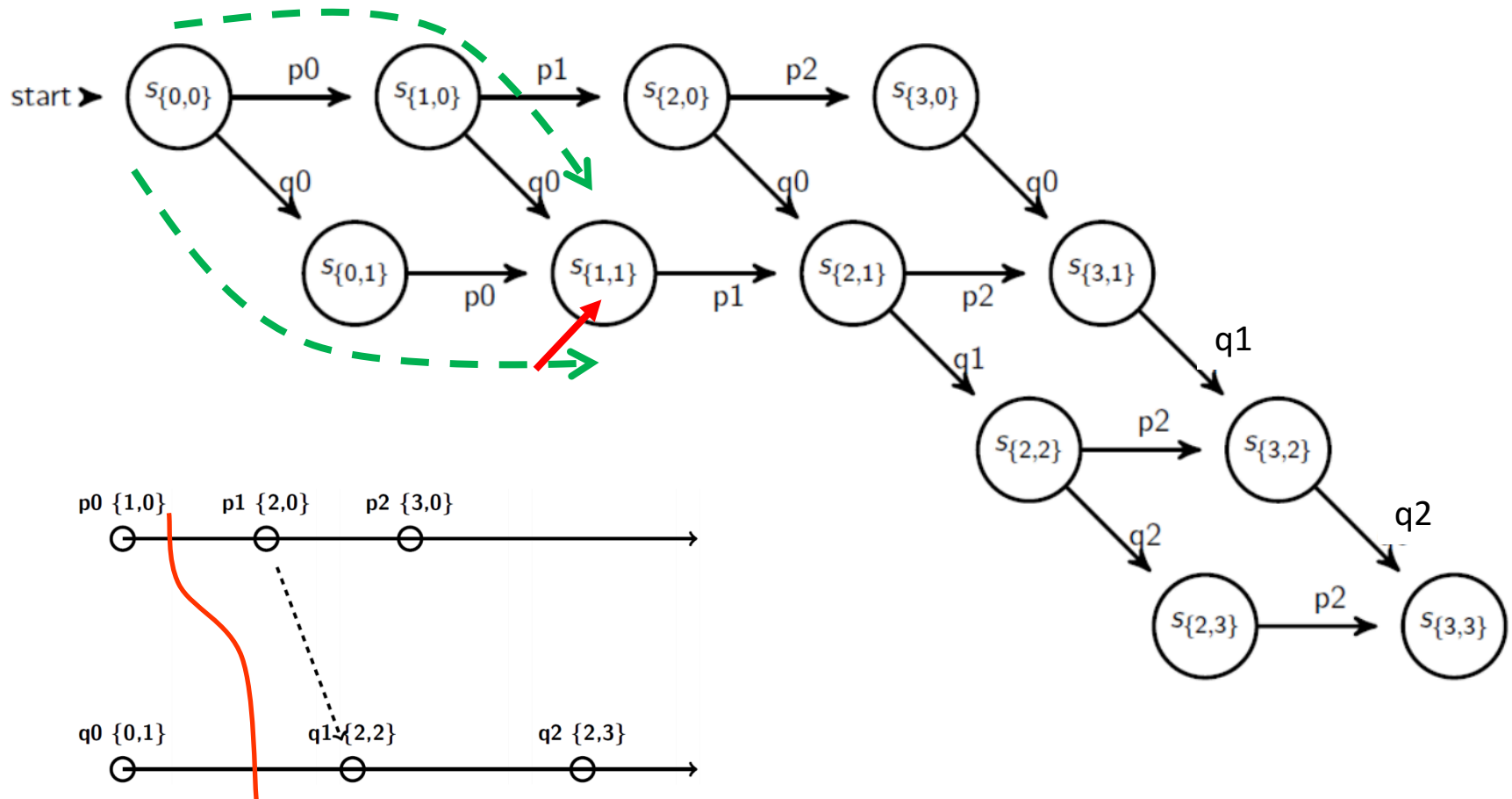


# State Transitions: Example

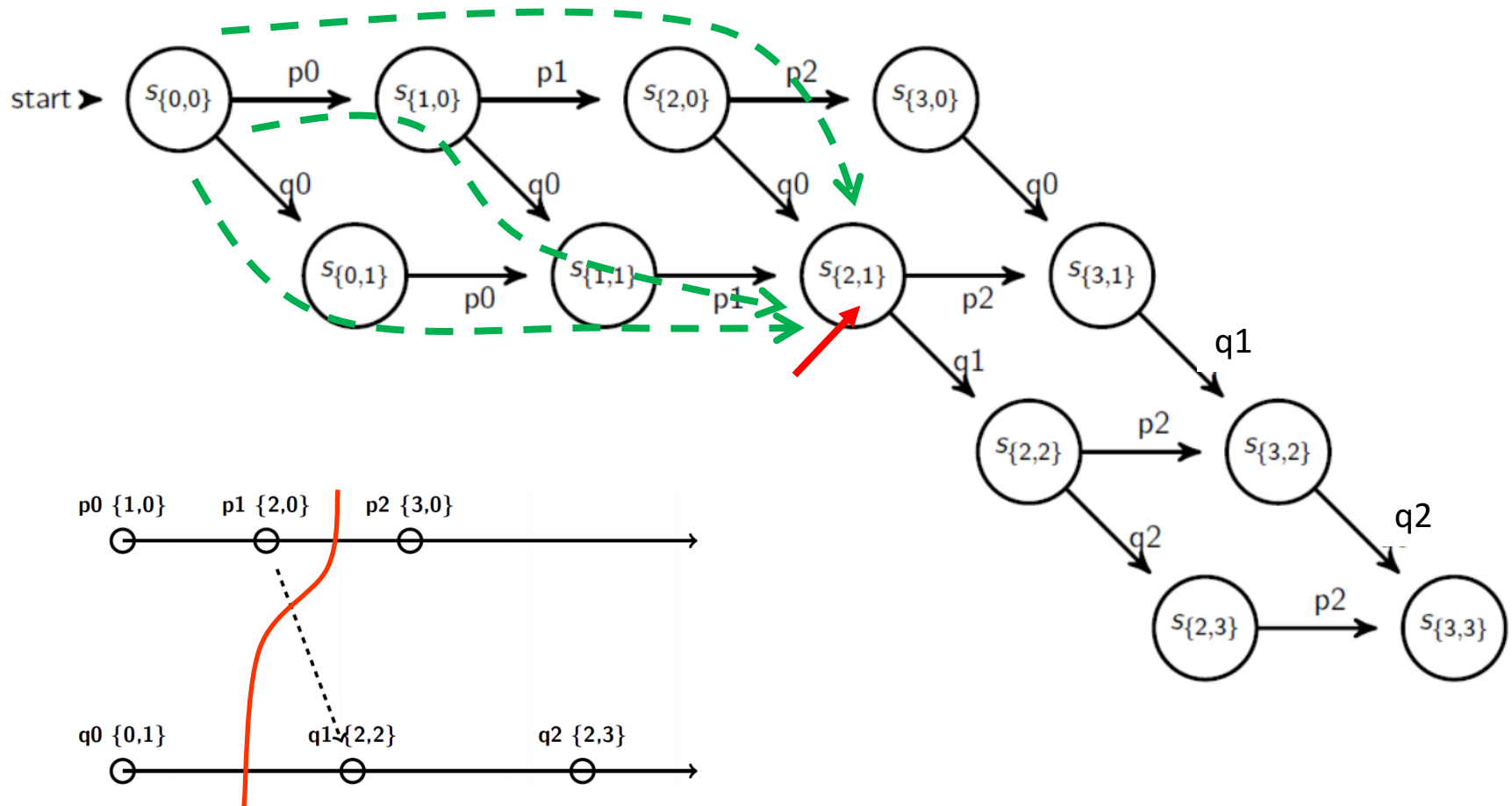
**Execution Lattice.** Each path represents a linearization.



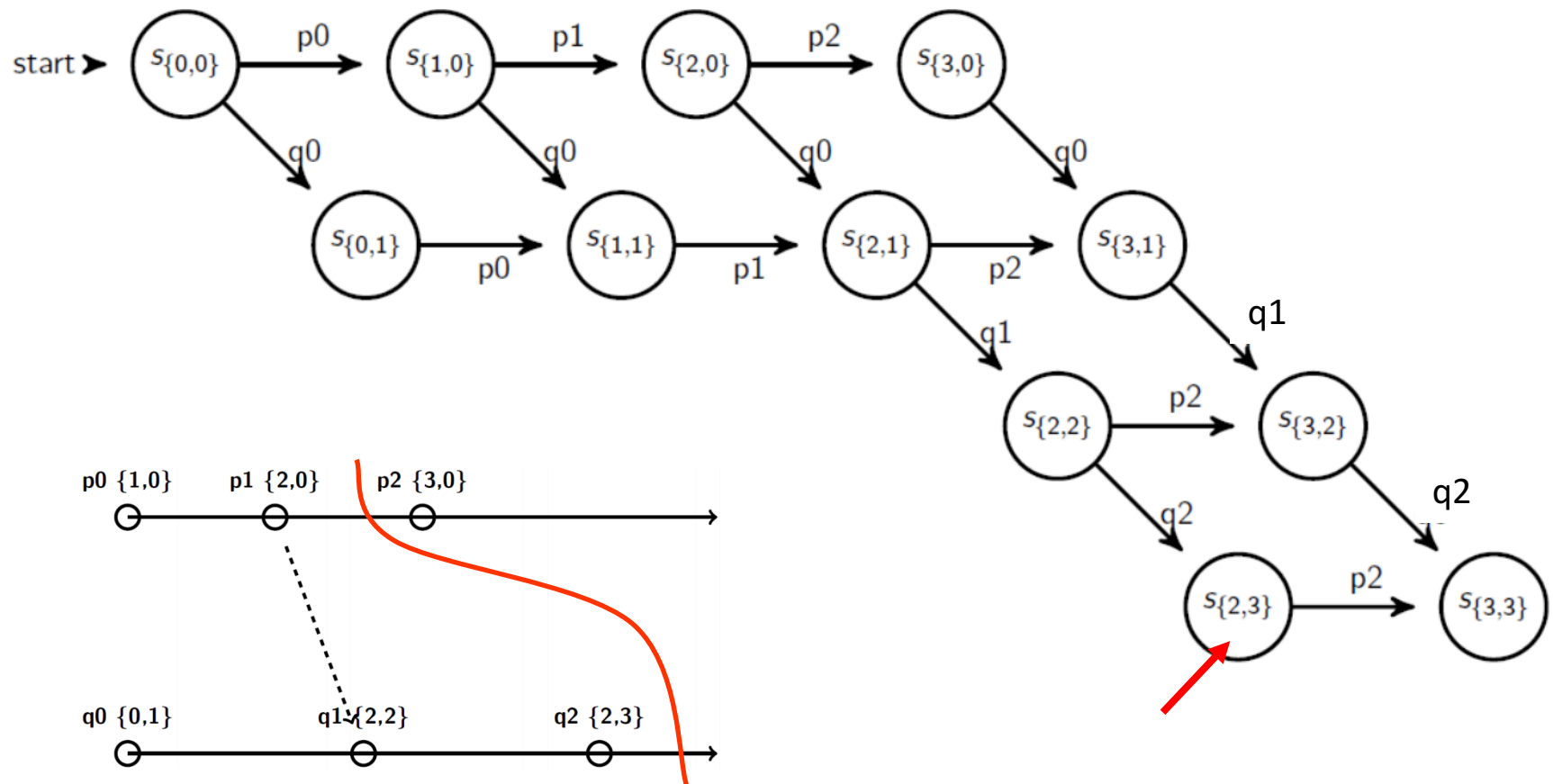
# State Transitions: Example



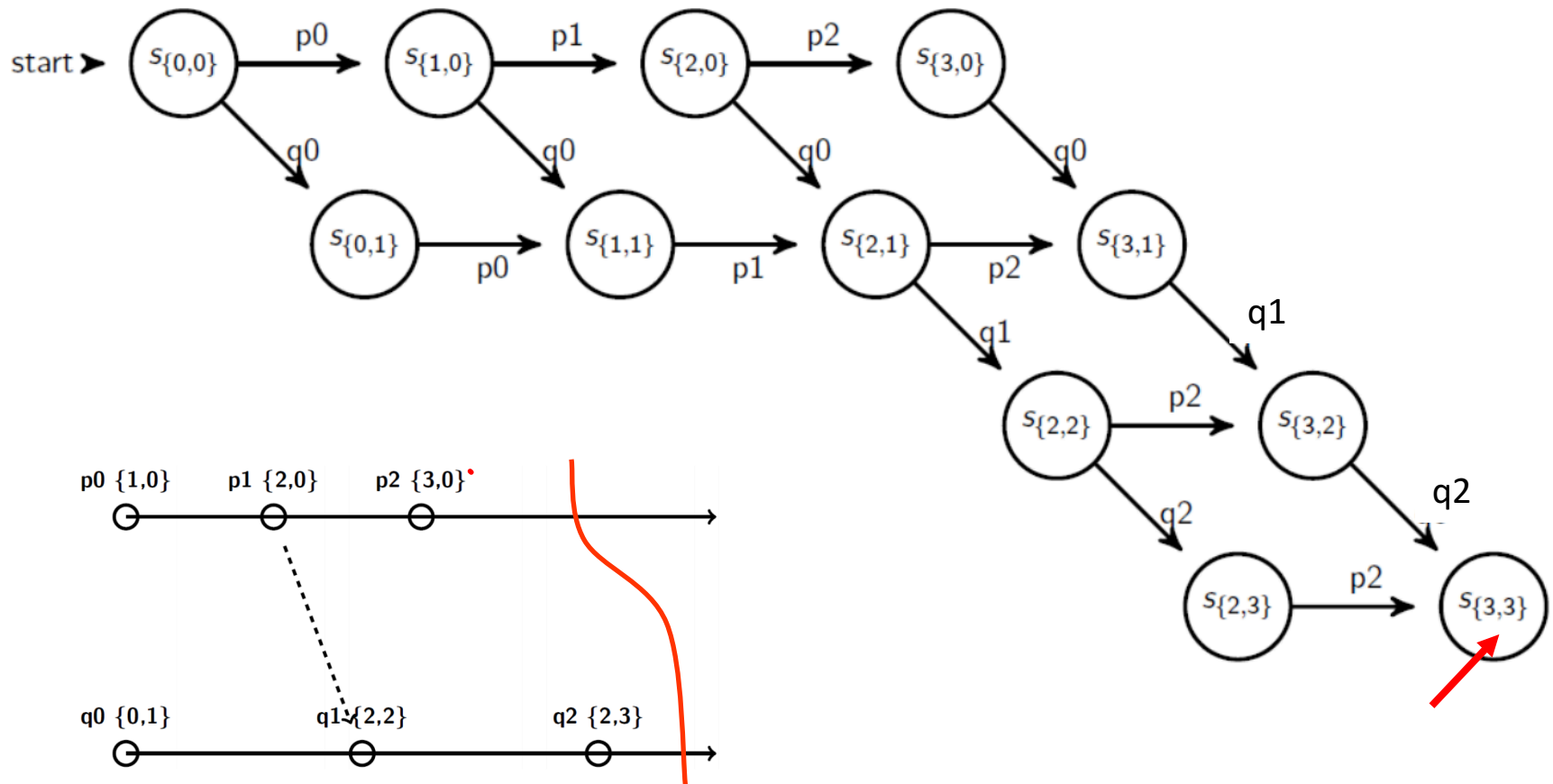
# State Transitions: Example



# State Transitions: Example



# State Transitions: Example



# More notations and definitions

- A **run** is a total ordering of events in  $H$  that is consistent with each  $h_i$ 's ordering.
- A **linearization** is a run consistent with happens-before ( $\rightarrow$ ) relation in  $H$ .
- Linearizations pass through consistent global states.
- A global state  $S_k$  is reachable from global state  $S_i$ , if there is a linearization that passes through  $S_i$  and then through  $S_k$ .
- The distributed system evolves as a series of transitions between global states  $S_0, S_1, \dots$



# Global State Predicates

- A global-state-predicate is a property that is *true* or *false* for a global state.
  - Is there a deadlock?
  - Has the distributed algorithm terminated?
- Two ways of reasoning about predicates (or system properties) as global state gets transformed by events.
  - Liveness
  - Safety

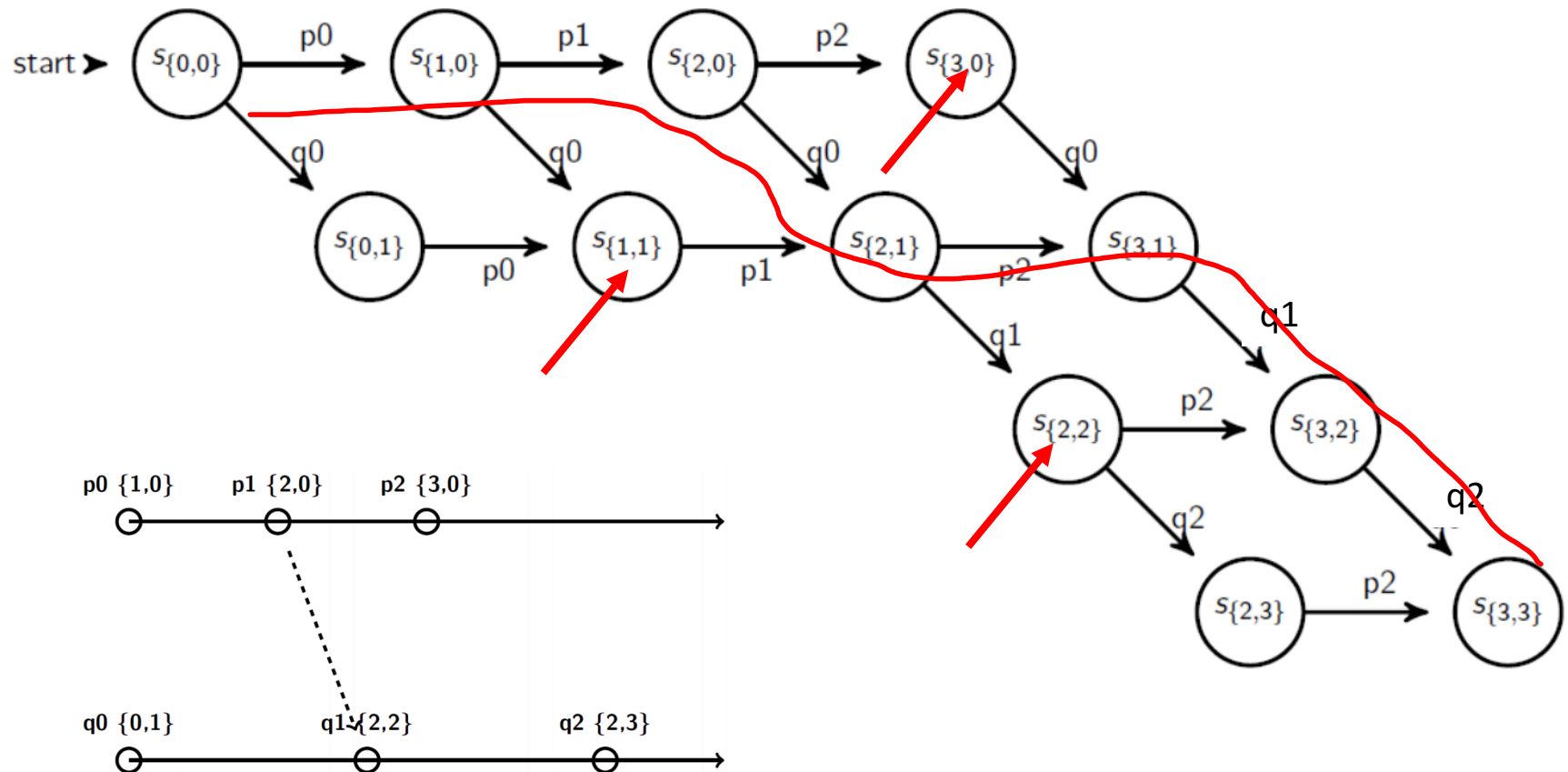
# Liveness

- **Liveness** = guarantee that something **good** will happen, **eventually**
- **Examples:**
  - A distributed computation will terminate.
  - “Completeness” in failure detectors: the failure will be detected.
  - All processes will eventually decide on a value.
- A global state  $S_0$  satisfies a **liveness** property  $P$  iff:
  - For all linearizations starting from  $S_0$ ,  $P$  is true for **some** state  $S_L$  reachable from  $S_0$ .
  - $\text{liveness}(P(S_0)) \equiv \forall L \in \text{linearizations from } S_0, L \text{ passes through a } S_L \text{ \& } P(S_L) = \text{true}$

# Liveness Example

If predicate is true only in the marked states, does it satisfy liveness?

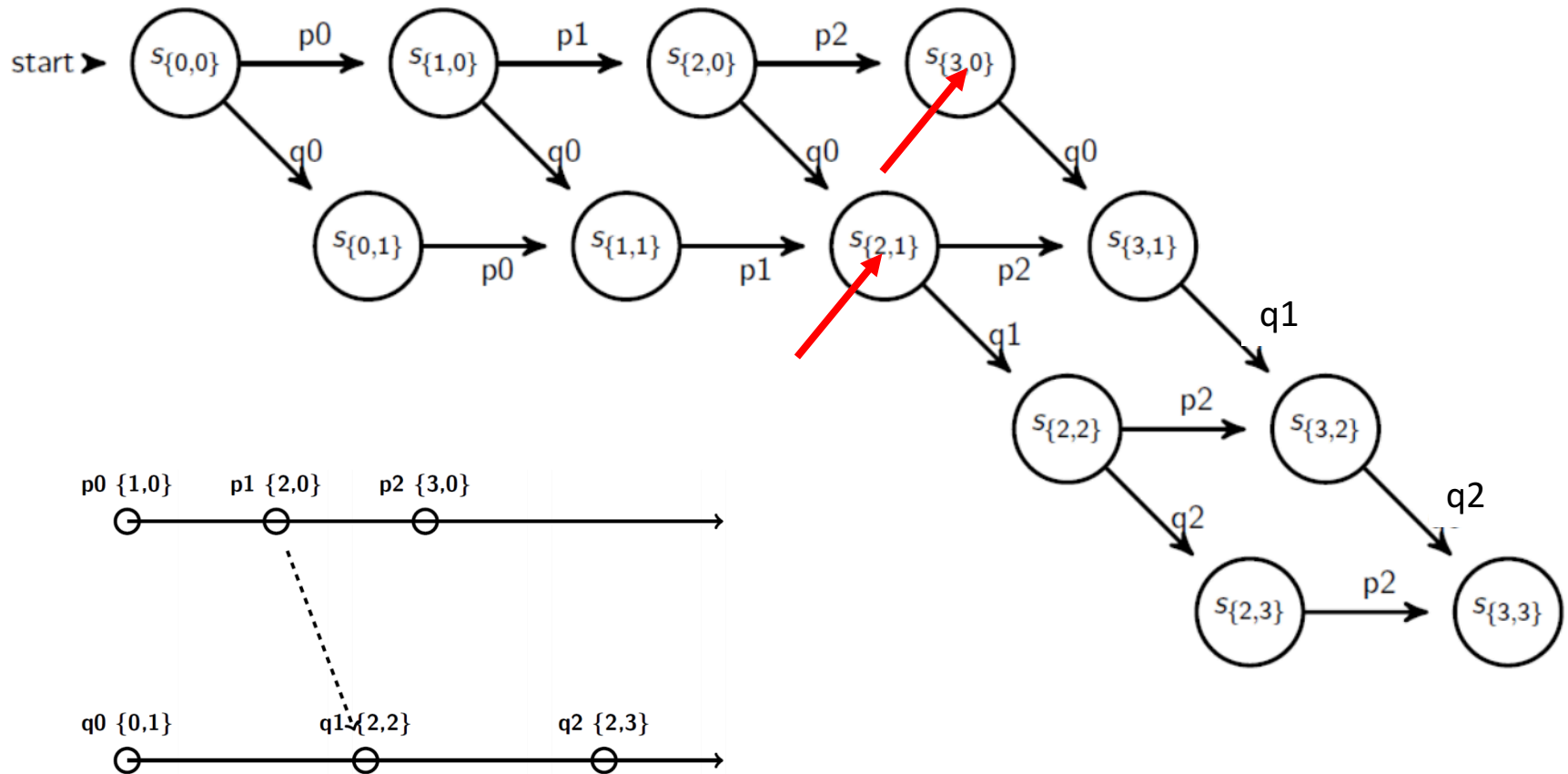
No



# Liveness Example

If predicate is true only in the marked states, does it satisfy liveness?

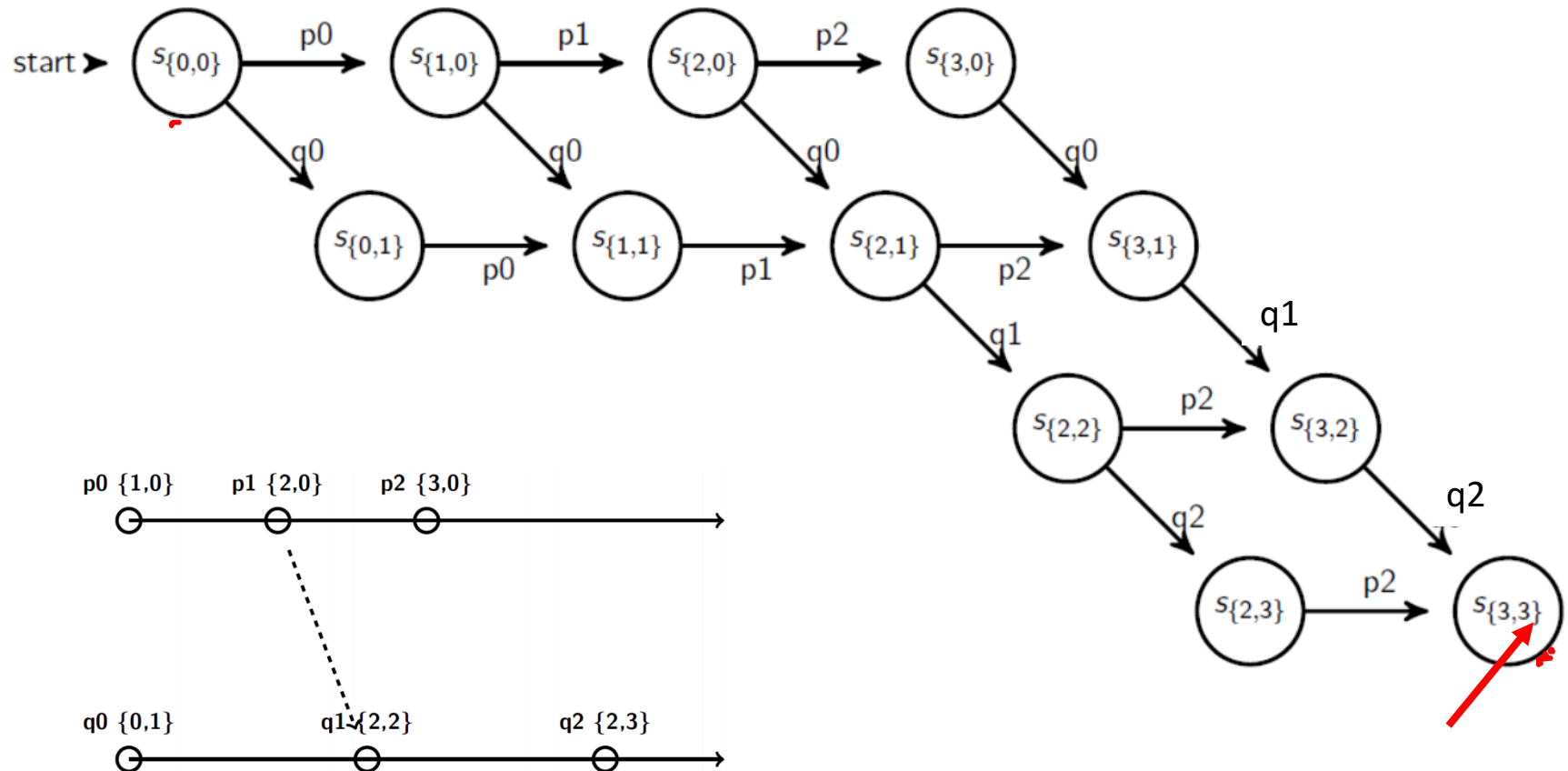
Yes



# Liveness Example

If predicate is true only in the marked states, does it satisfy liveness?

Yes



# Liveness

- **Liveness** = guarantee that something **good** will happen, **eventually**
- **Examples:**
  - A distributed computation will terminate.
  - “Completeness” in failure detectors: the failure will be detected.
  - All processes will eventually decide on a value.
- A global state  $S_0$  satisfies a **liveness** property  $P$  iff:
  - $\text{liveness}(P(S_0)) \equiv \forall L \in \text{linearizations from } S_0, L \text{ passes through a } S_L \text{ \& } P(S_L) = \text{true}$
  - For any linearization starting from  $S_0$ ,  $P$  is true for **some** state  $S_L$  reachable from  $S_0$ .

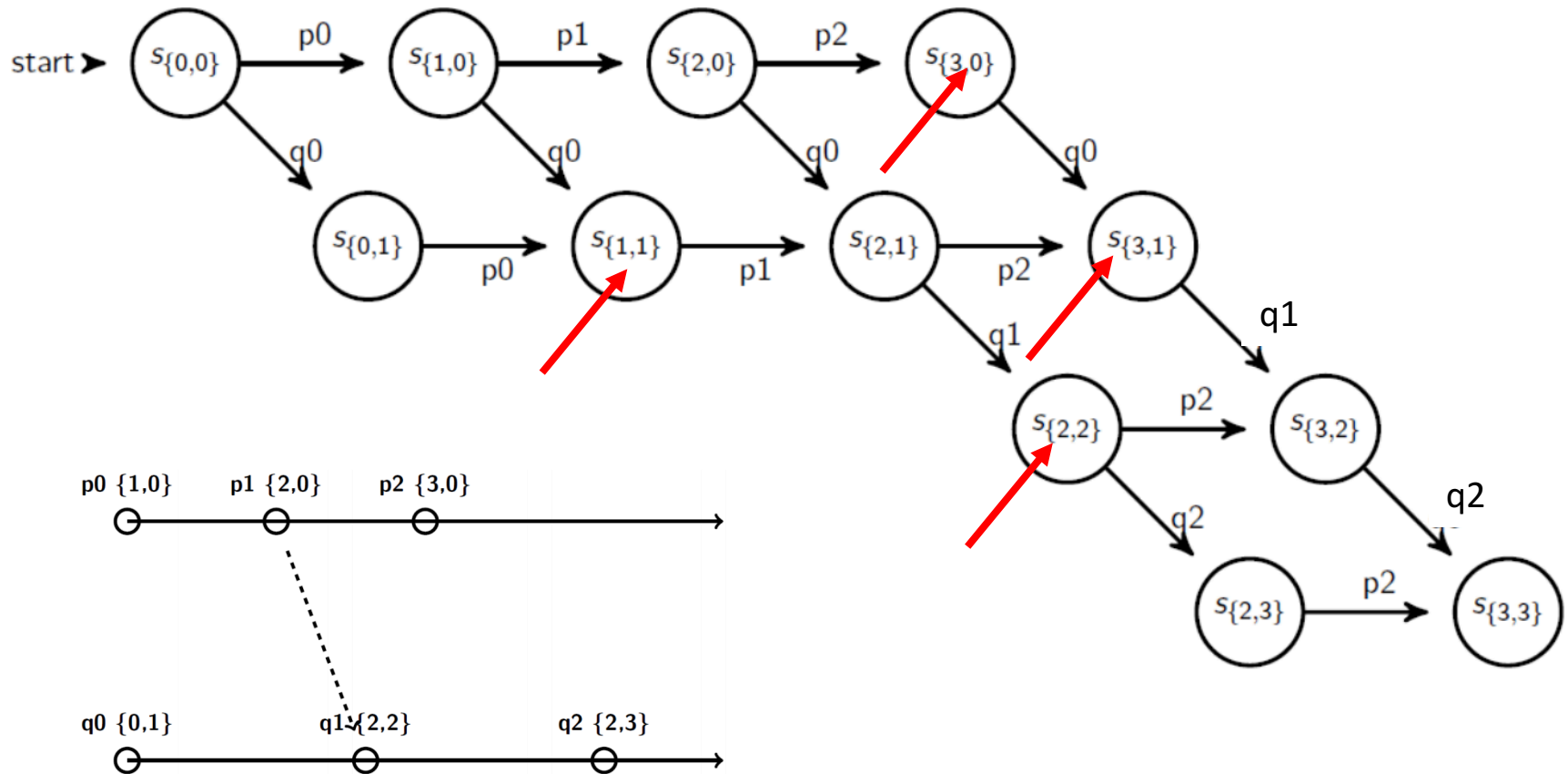
# Safety

- **Safety** = guarantee that something **bad** will **never** happen.
- **Examples:**
  - There is no deadlock in a distributed transaction system.
  - “Accuracy” in failure detectors: an alive process is not detected as failed.
  - No two processes decide on different values.
- A global state  $S_0$  satisfies a **safety** property  $P$  iff:
  - For **all** states  $S$  reachable from  $S_0$ ,  $P(S)$  is true.
  - $\text{safety}(P(S_0)) \equiv \forall S \text{ reachable from } S_0, P(S) = \text{true}.$

# Safety Example

If predicate is true only in the marked states, does it satisfy safety?

No

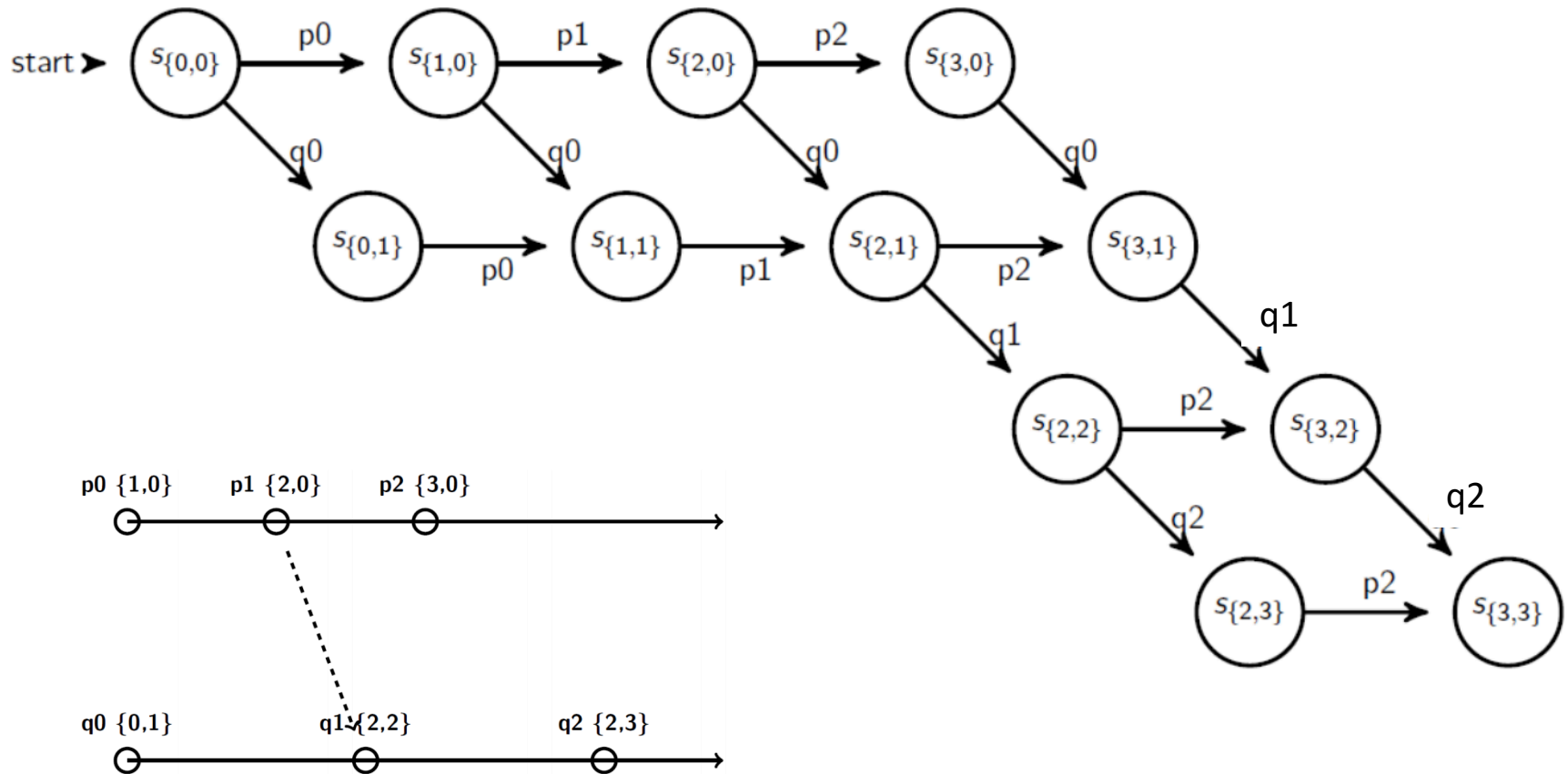




# Safety Example

If predicate is true only in the **unmarked** states, does it satisfy safety?

Yes

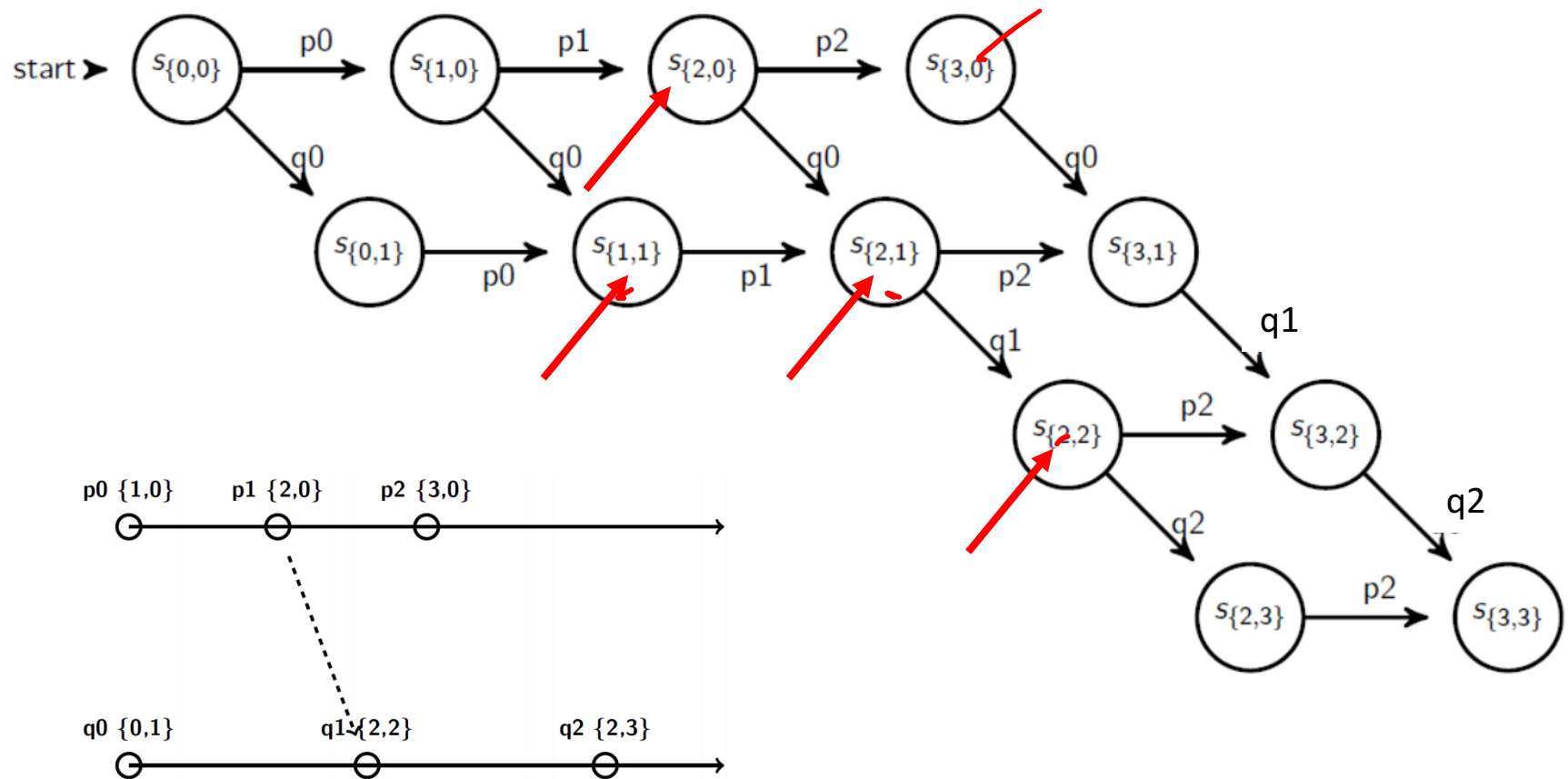


# Safety

- **Safety** = guarantee that something **bad** will **never** happen.
- **Examples:**
  - There is no deadlock in a distributed transaction system.
  - “Accuracy” in failure detectors: an alive process is not detected as failed.
  - No two processes decide on different values.
- A global state  $S_0$  satisfies a **safety** property  $P$  iff:
  - $\text{safety}(P(S_0)) \equiv \forall S \text{ reachable from } S_0, P(S) = \text{true}.$
  - For **all** states  $S$  reachable from  $S_0$ ,  $P(S)$  is true.

# Liveness Example

Technically satisfies liveness, but difficult to capture or reason about.



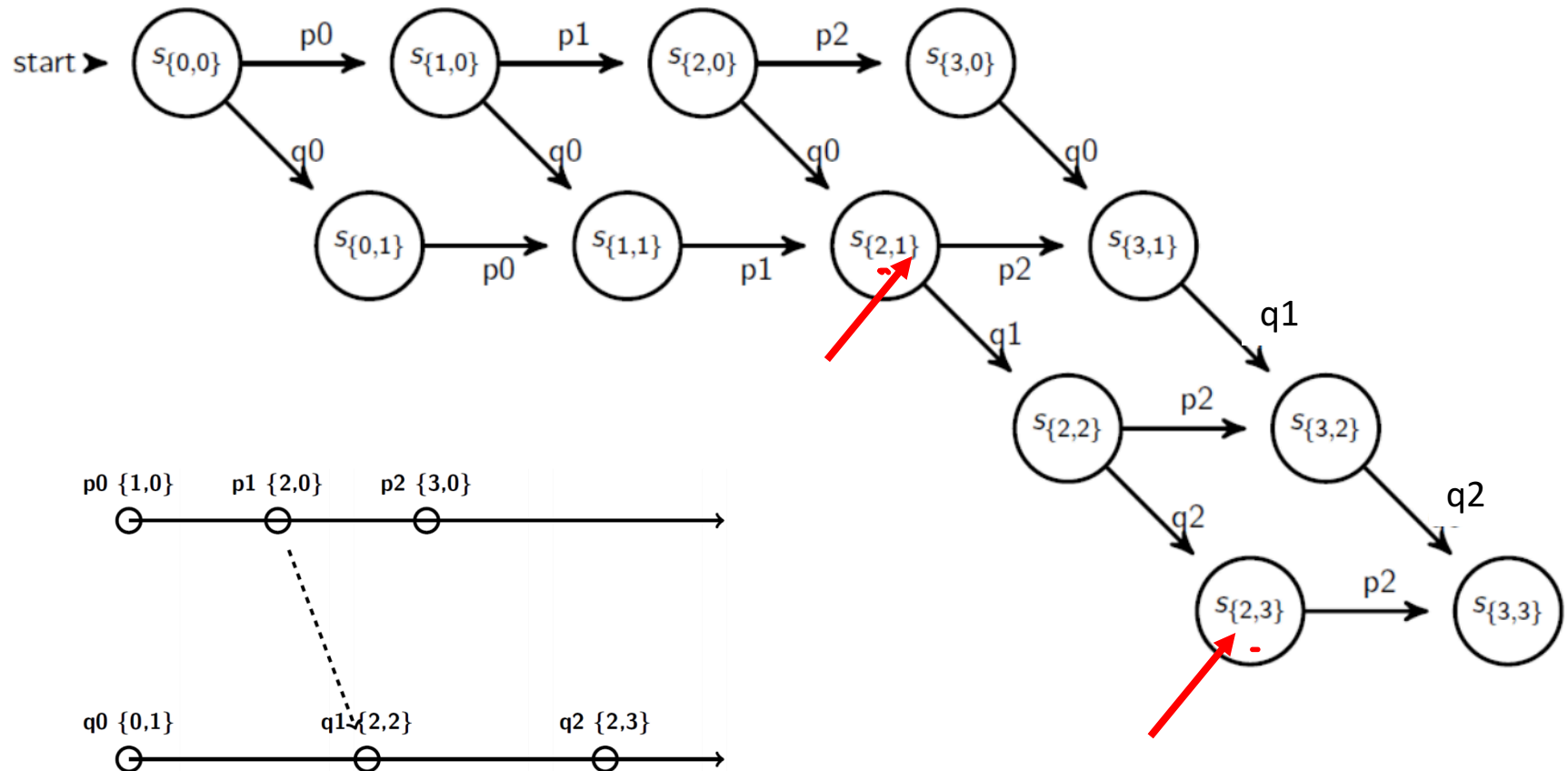
# Stable Global Predicates

- once true, stays true forever afterwards (for stable liveness)

# Stable Global Predicates

If predicate is true only in the marked states, is it true in a stable way?

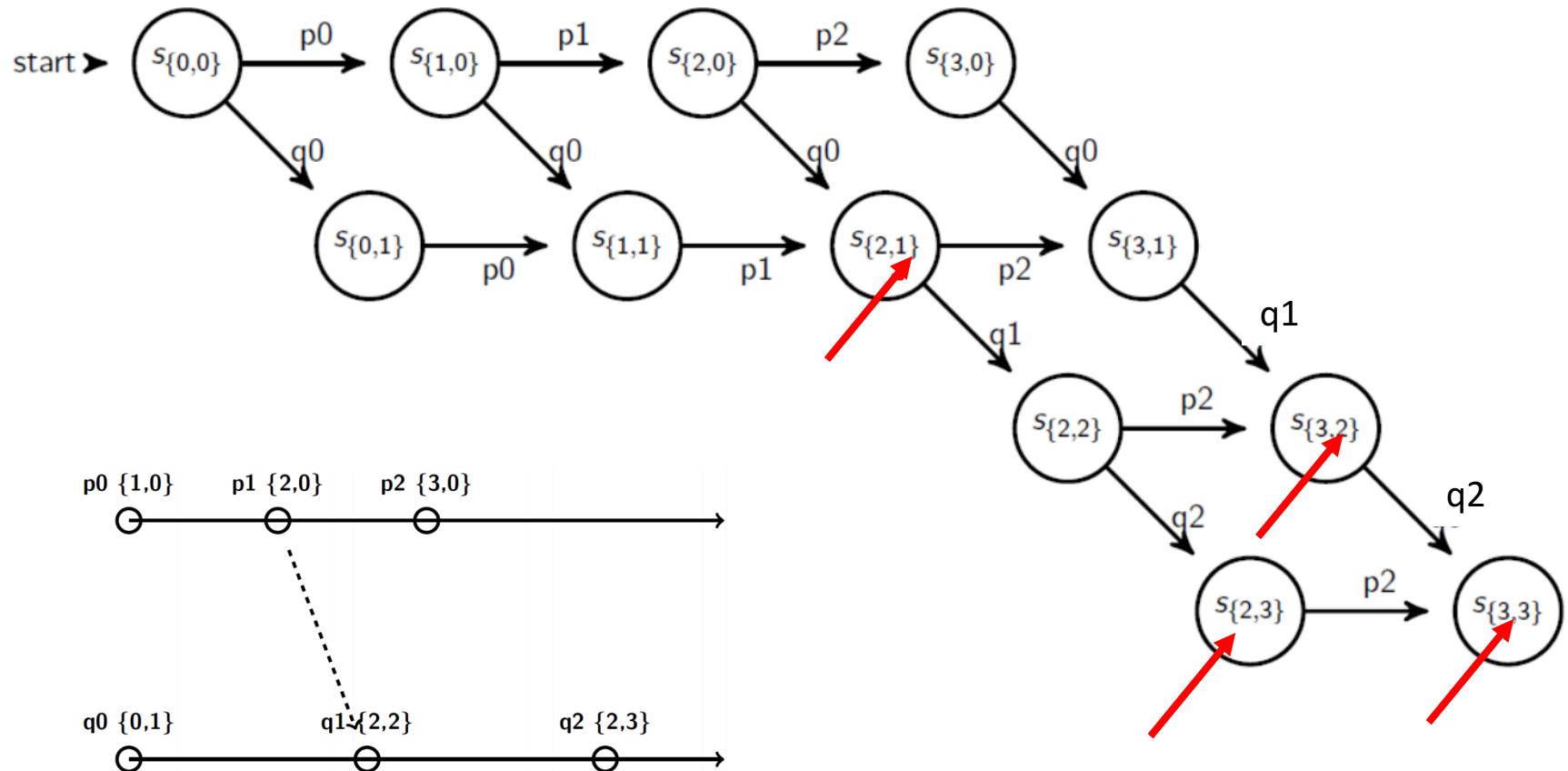
No



# Stable Global Predicates

If predicate is true only in the marked states, is it true in a stable way?

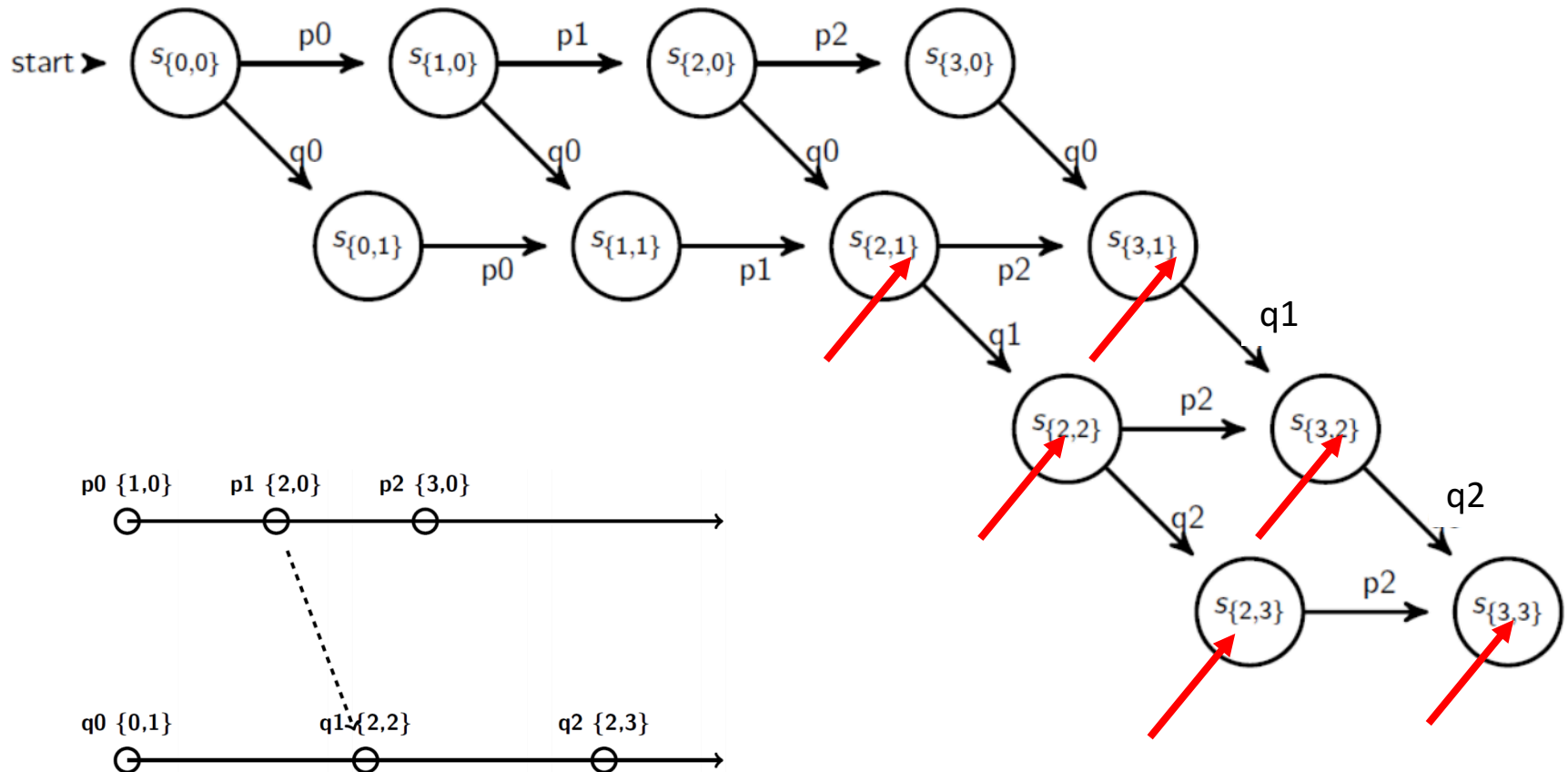
No



# Stable Global Predicates

If predicate is true only in the marked states, is it true in a stable way?

Yes



# Stable Global Predicates

- once true for a state  $S$ , stays true for all states reachable from  $S$  (for stable liveness)
- once false for a state  $S$ , stays false for all states reachable from  $S$  (for stable non-safety)
- Stable liveness examples (once true, always true)
  - Computation has terminated.
- Stable non-safety examples (once false, always false)
  - There is no deadlock.
  - An object is not orphaned.
- *All stable global properties can be detected using the Chandy-Lamport algorithm.*



# Global Snapshot Summary

- The ability to calculate global snapshots in a distributed system is very important.
- But don't want to interrupt running distributed application.
- Chandy-Lamport algorithm calculates global snapshot.
- Obeys causality (creates a consistent cut).
- Can be used to detect global properties.
- Safety vs. Liveness.

# Today's agenda

- Global State (contd.)

- Chapter 14.5

- **Multicast**

- Chapter 15.4

- **Goal:** reason about desirable properties for message delivery among a group of processes.

# Communication modes

- **Unicast**
  - Messages are sent from exactly one process to one process.
- **Broadcast**
  - Messages are sent from exactly one process to all processes on the network.
- **Multicast**
  - Messages broadcast within a group of processes.
  - A multicast message is sent from any one process to a group of processes on the network.

# Where is multicast used?

- Distributed storage
  - Write to an object are multicast across replica servers.
  - Membership information (e.g., heartbeats) is multicast across all servers in cluster.
- Online scoreboards (ESPN, French Open, FIFA World Cup)
  - Multicast to group of clients interested in the scores.
- Stock Exchanges
  - Group is the set of broker computers.
- .....

# Communication modes

- **Unicast**

- Messages are sent from exactly one process to one process.
  - *Best effort*: if a message is delivered it would be intact; no reliability guarantees.
  - *Reliable*: guarantees delivery of messages.
  - *In order*: messages will be delivered in the same order that they are sent.

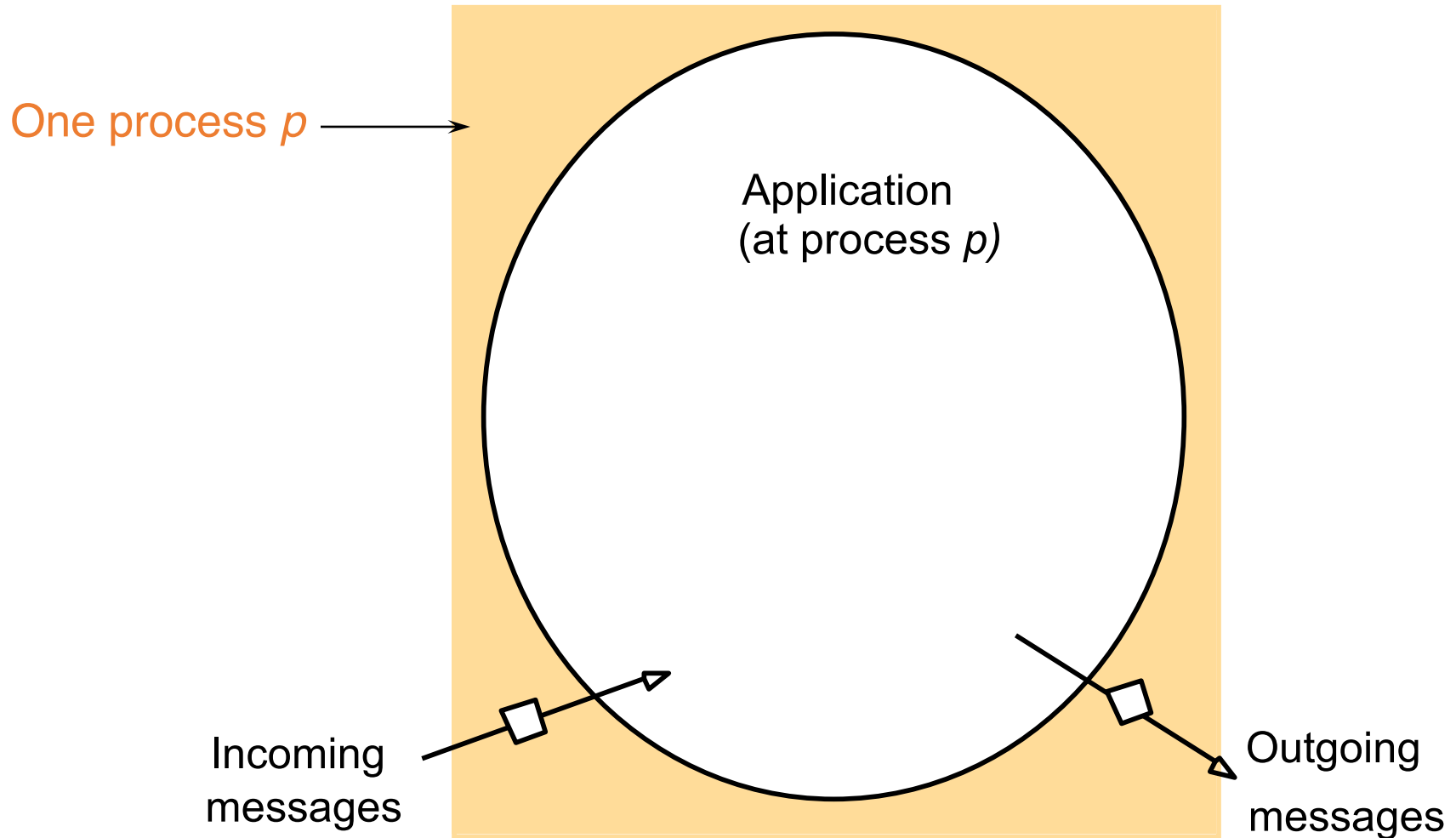
- **Broadcast**

- Messages are sent from exactly one process to all processes on the network.

- **Multicast**

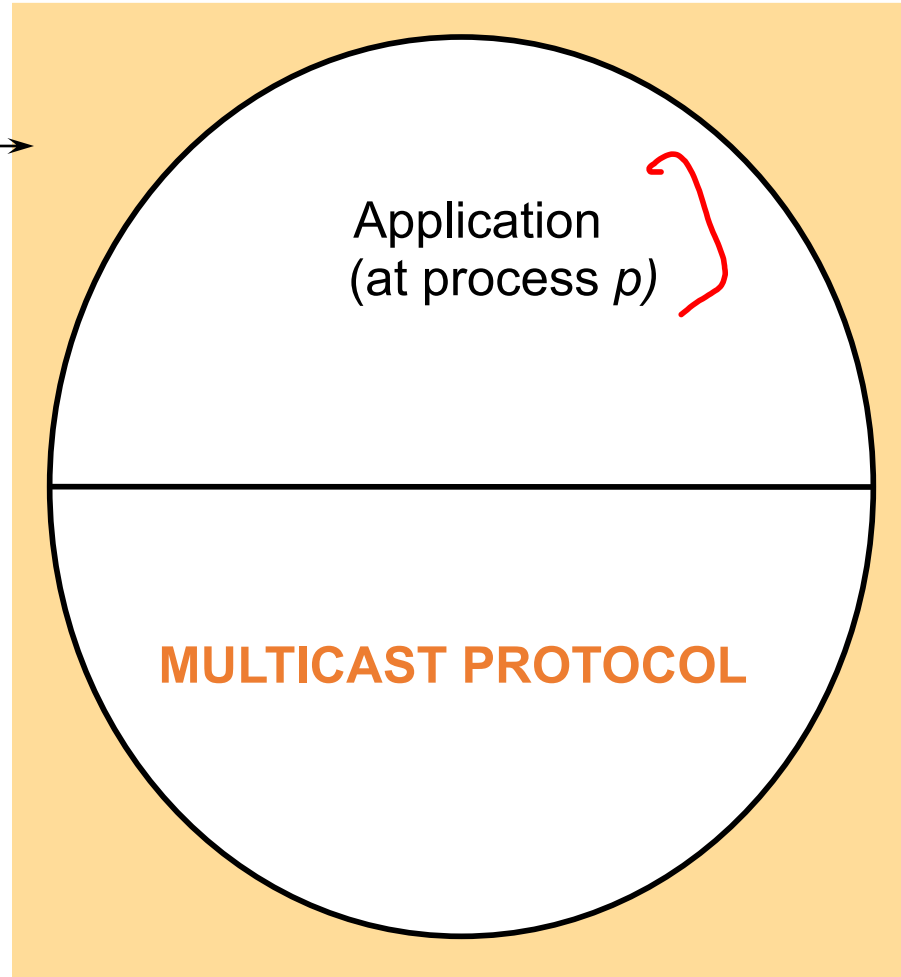
- Messages broadcast within a group of processes.
- A multicast message is sent from any one process to the group of processes on the network.
- *How do we define (and achieve) reliable or ordered multicast?*

# What we are designing in this class?

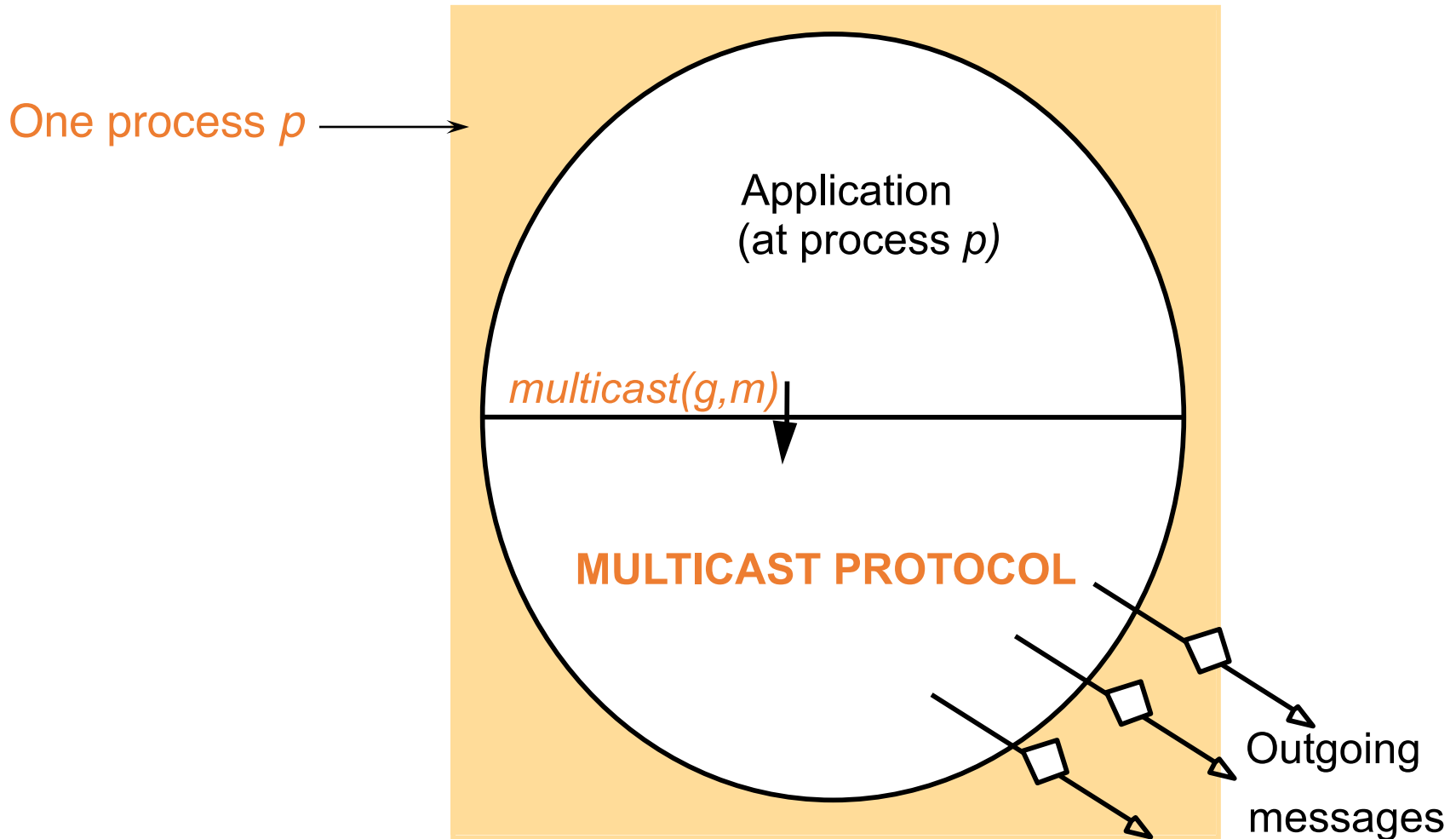


# What we are designing in this class?

One process  $p$



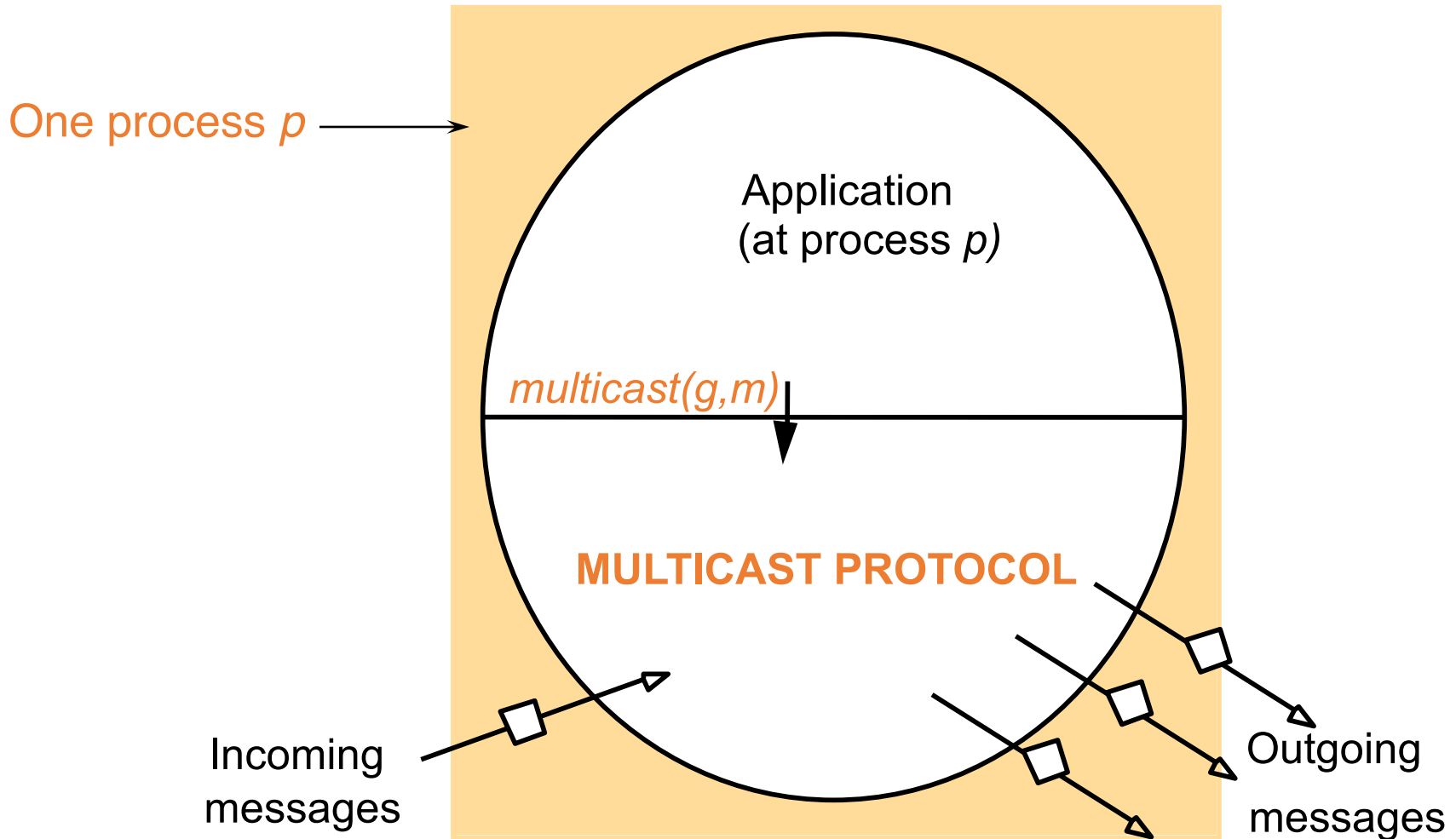
# What we are designing in this class?



' $g$ ' is a multicast group that also includes the process ' $p$ '.

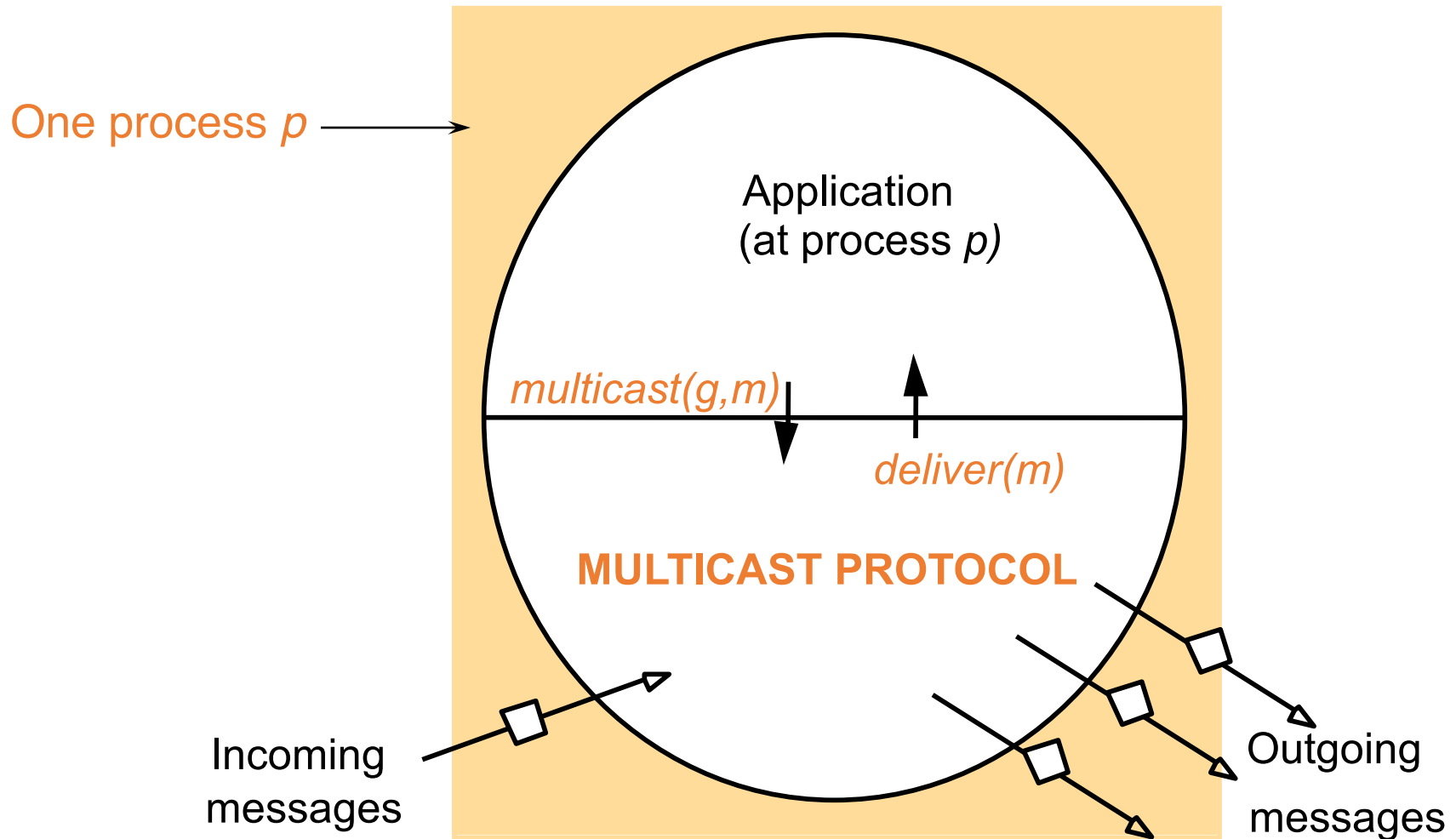


# What we are designing in this class?



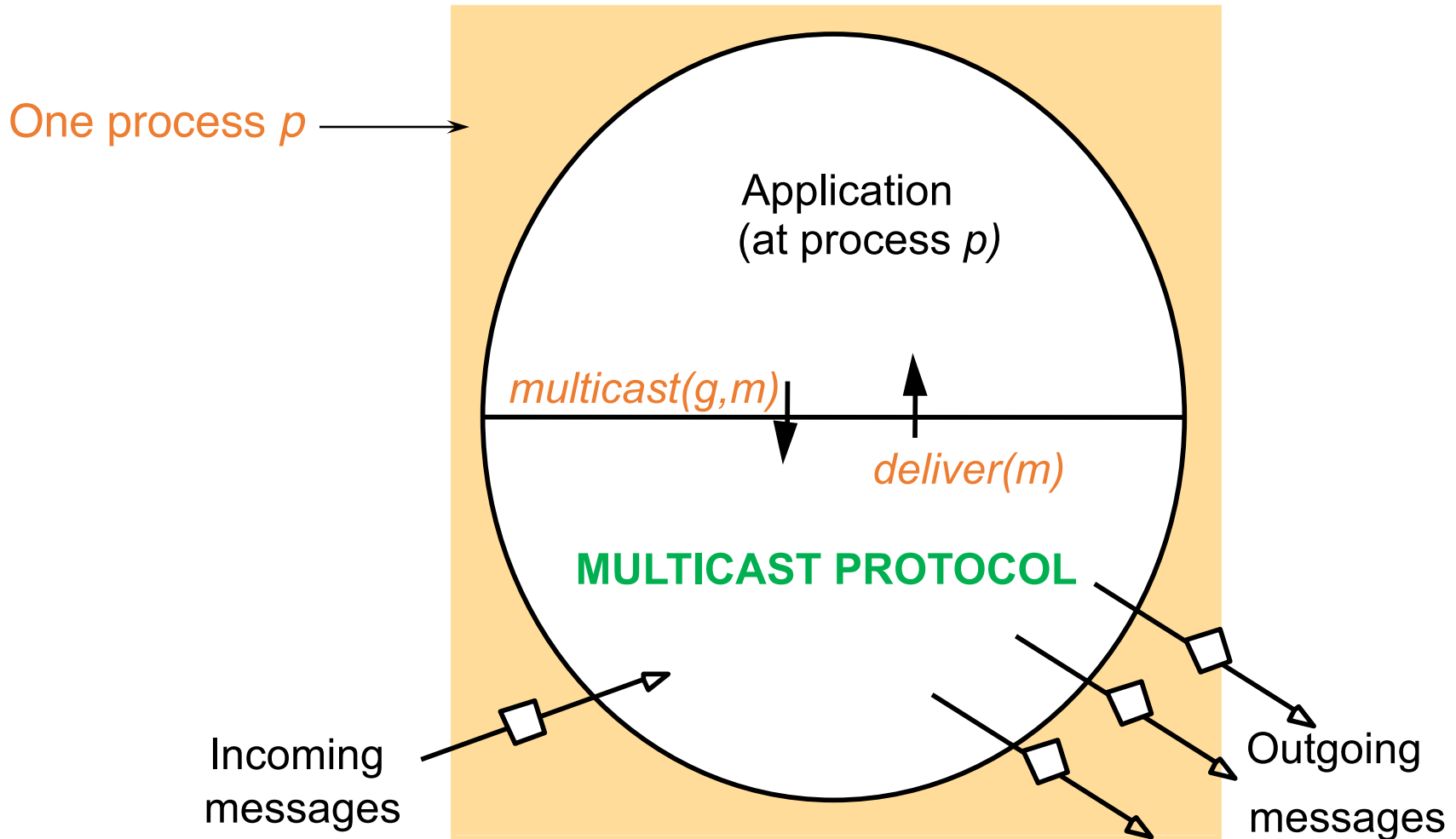
' $g$ ' is a multicast group that also includes the process ' $p$ '.

# What we are designing in this class?



' $g$ ' is a multicast group that also includes the process ' $p$ '.

# What we are designing in this class?



' $g$ ' is a multicast group that also includes the process ' $p$ '.

# Basic Multicast (B-Multicast)

- Straightforward way to implement B-multicast:
  - use a reliable one-to-one send (unicast) operation:  
B-multicast(group  $g$ , message  $m$ ):  
for each process  $p$  in  $g$ , send ( $p,m$ ).  
receive( $m$ ): B-deliver( $m$ ) at  $p$ .
- Guarantees: message is eventually delivered to the group if:
  - Processes are non-faulty.
  - The unicast “send” is reliable.
  - *Sender does not crash.*
- *Can we provide reliable delivery even after sender crashes?*
  - *What does this mean?*

# Reliable Multicast (R-Multicast)

- **Integrity:** A *correct* (i.e., non-faulty) process  $p$  delivers a message  $m$  at most once.
  - *Assumption: no process sends **exactly** the same message twice*
- **Validity:** If a *correct* process multicasts (sends) message  $m$ , then it will *eventually* deliver  $m$  to itself.
  - *Liveness for the sender.*
- **Agreement:** If a *correct* process delivers message  $m$ , then all the other *correct* processes in  $\text{group}(m)$  will *eventually* deliver  $m$ .
  - *All or nothing.*
- Validity and agreement together ensure overall liveness: if some correct process multicasts a message  $m$ , then, all correct processes deliver  $m$  too.

How to achieve R-multicast? To be continued in next class....